

Continuous Latent Variables

Prof. Nicholas Zabaras

Center for Informatics and Computational Science

<https://cics.nd.edu/>

*University of Notre Dame
Notre Dame, Indiana, USA*

Email: nzabaras@gmail.com

URL: <https://www.zabaras.com/>

March 21, 2019



Contents - PCA

- Continuous Latent Variable Model, Generative Point of View, Factor Analysis, Unidentifiability, Mixture of Factor Analysers
- Principal Component Analysis, Maximum variance formulation, Minimum-error formulation, PCA Reconstruction, PCA and SVD
- Applications of PCA, Off-line Digit Images, Whitening of the data with PCA, Old Faithful Data Set, Oil Flow Data Set, PCA for Visualization, PCA Vs. Fisher Discriminant
- PCA for High-Dimensional Data

Following closely Chris Bishop's PRML book (Chapter 12)



Contents - Probabilistic PCA

- Probabilistic PCA, Generative View Point, Predictive Distribution,
Posterior Distribution, Maximum likelihood PCA, Efficient Evaluation of PPCA, Maximum Likelihood PCA Vs. Standard PCA, Probabilistic PCA – Degrees of Freedom

- EM algorithm for PCA, EM Cost, EM Missing Values, EM for $\sigma^2 \rightarrow 0$,
EM Examples, EM Algorithm and Standard PCA Revisited

Following closely Chris Bishops' PRML book (Chapter 12)



Contents-Nonlinear Latent Variable Models

- Independent component analysis, Fast ICA Algorithm, Modeling the Source Densities, EM Implementation, Other Estimation Techniques, Maximizing non-Gaussianity, Maximizing Mutual Information
- Autoassociative neural networks
- Modelling nonlinear manifolds, Mixture of PCA, Principal Curves, Multidimensional Scaling, Locally Linear Embedding, ISOMAP, Latent Trait Models, Density Network
- Generative Topographic Mapping

Following closely Chris Bishop's PRML book (Chapter 12)

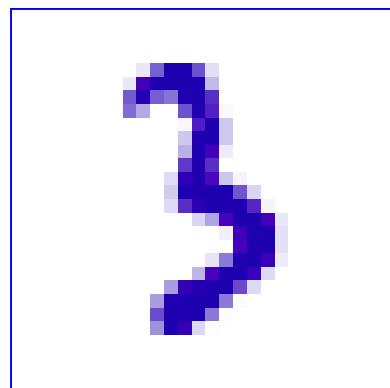


Continuous Latent Variables - Introduction



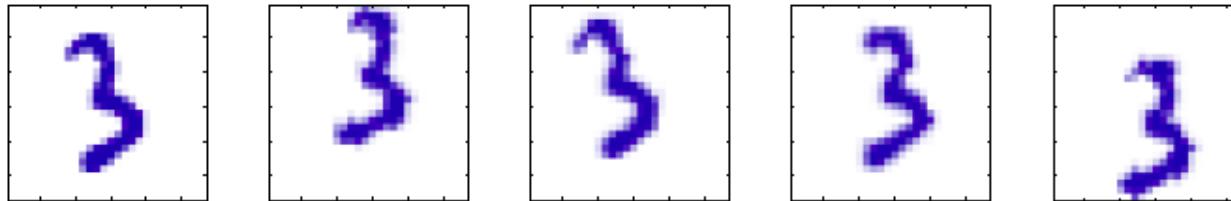
Continuous Latent Variables

- In many data sets, the data points lie close to a manifold of much lower dimensionality than that of the original data space.
- Consider a data set constructed by taking one of the off-line digits, represented by a 64×64 pixel grey-level image.
- Embed this in a larger image of size 100×100 by padding with pixels having the value zero (white pixels).
- Each image is represented by a point in 10,000 dimensional space.



Continuous Latent Variables

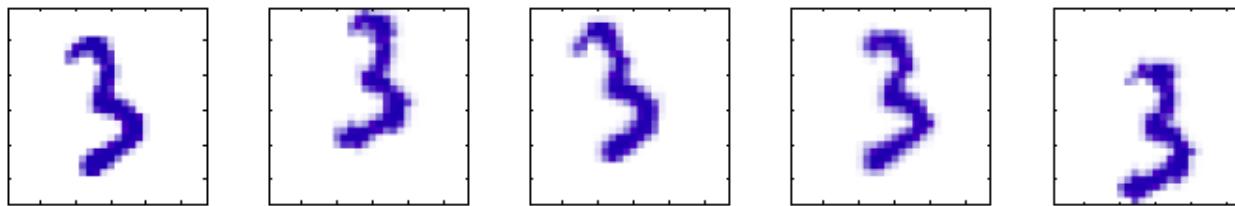
- We create multiple copies in which the location and orientation of the digit is varied at random.



- In this data set, there are 3 DOF of variability:
 - Vertical displacement
 - Horizontal displacement and
 - Rotation

Intrinsic Dimensionality of a Data Set

- The intrinsic dimensionality of the data set is three (vertical and horizontal displacement and rotation)



- The manifold is non-linear: when translating the digit past a particular pixel, that pixel value goes from 0 (white) to 1 (black) and back to zero again. This is a nonlinear function of the digit position.
- The translation and rotation parameters are latent variables: we observe only the image vectors without knowing the translation or rotation variables used to create them.

Other Latent Variables

- For real digit image data, additional DOF arise from
 - Scaling
 - Complex deformations due to the variability of an individual's writing style,
 - Etc.
- The number of such degrees of freedom will still be small compared to the dimensionality of the data set.
- For data compression, we are interested to explore the manifold structure.

Generative Point of View



Generative Point of View

- The data points are often **not** confined precisely to a smooth low-dimensional manifold.
- The departure of data points from the manifold is interpreted as ‘noise’.
- **Generative view:**
 - (a) Select a point within the manifold according to some latent variable distribution $p(z)$ and
 - (b) Generate an observed data point by adding noise, drawn from some conditional distribution of the data variables given the latent variables $p(x|z)$.

Generative Point of View

□ Example – Linear Gaussian latent variable model:

- (a) assume Gaussian distributions for the latent variables
 $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I})$ and the observed variables
 $p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}|W\mathbf{z} + \boldsymbol{\mu}, \sigma^2\mathbf{I})$
 - (b) make use of a linear-Gaussian dependence of the observed variables on the state of the latent variables,
 $p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}|W\mathbf{z} + \boldsymbol{\mu}, \sigma^2\mathbf{I})$
- This leads to a probabilistic formulation of principal component analysis (PCA).

Generative Point of View

□ Example – Factor Analysis:

- (a) assume Gaussian distributions for the latent variables

$$p(\mathbf{z}_i) = \mathcal{N}(\mathbf{z}_i | \boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$$

and the observed variables

$$p(\mathbf{x}_i | \mathbf{z}_i, \boldsymbol{\theta}) = \mathcal{N}(\mathbf{x}_i | \mathbf{W}\mathbf{z}_i + \boldsymbol{\mu}, \boldsymbol{\Psi}),$$

$$\mathbf{x}_i \in \mathbb{R}^D, \mathbf{z}_i \in \mathbb{R}^M, \mathbf{W} \in \mathbb{R}^{D \times M}, \boldsymbol{\Psi} \in \mathbb{R}^{D \times D}$$

- (b) We take $\boldsymbol{\Psi}$ to be diagonal. This overall model is called **factor analysis** or **FA**.
- The special case in which $\boldsymbol{\Psi} = \sigma^2 \mathbf{I}$ is called **probabilistic principal components analysis** or **PPCA**.

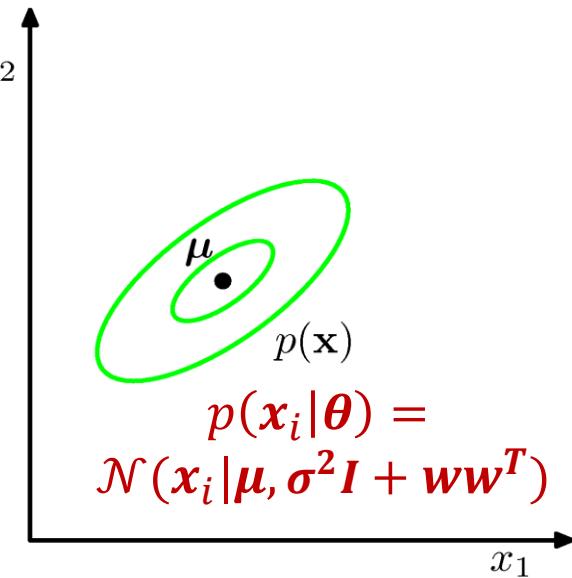
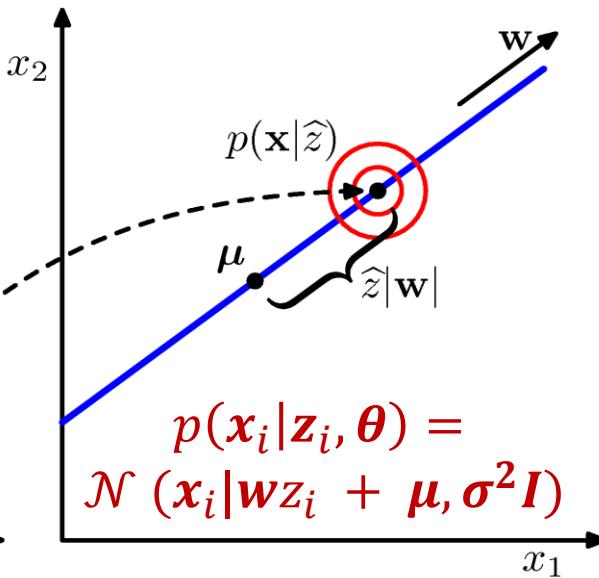
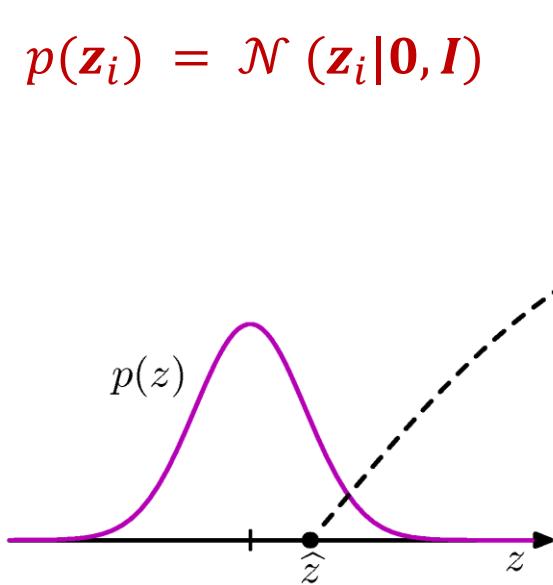


Generative Point of View

- PPCA is illustrated for $D = 2, M = 1$ (latent dimension).

$$p(\mathbf{x}_i|\boldsymbol{\theta}) = \int \mathcal{N}(\mathbf{x}_i|W\mathbf{z}_i + \boldsymbol{\mu}, \boldsymbol{\Psi}) \mathcal{N}(\mathbf{z}_i|\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0) d\mathbf{z}_i = \mathcal{N}(\mathbf{x}_i|W\boldsymbol{\mu}_0 + \boldsymbol{\mu}, \boldsymbol{\Psi} + W\boldsymbol{\Sigma}_0 W^T)$$

$$p(\mathbf{z}_i) = \mathcal{N}(\mathbf{z}_i|\mathbf{0}, \mathbf{I})$$



$$p(\mathbf{z}_i) = \mathcal{N}(\mathbf{z}_i|\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$$

$$p(\mathbf{x}_i|\mathbf{z}_i, \boldsymbol{\theta}) = \mathcal{N}(\mathbf{x}_i|W\mathbf{z}_i + \boldsymbol{\mu}, \boldsymbol{\Psi})$$

$$p(\mathbf{x}_i|\boldsymbol{\theta}) =$$

$$\mathcal{N}(\mathbf{x}_i|\boldsymbol{\mu}, \sigma^2 \mathbf{I} + \mathbf{w}\mathbf{w}^T)$$

- Based on the obtained marginal, we see that without sacrificing generality, we can simplify as:

$$p(\mathbf{z}_i) = \mathcal{N}(\mathbf{z}_i|\mathbf{0}, \mathbf{I})$$

$$p(\mathbf{x}_i|\boldsymbol{\theta}) = \mathcal{N}(\mathbf{x}_i|\boldsymbol{\mu}, \boldsymbol{\Psi} + \mathbf{W}\mathbf{W}^T)$$

Factor Analysis



Factor Analysis

- FA is a low rank parameterization of a MVN.
- FA can be thought of as a way of specifying a joint density model on x using a small number of parameters

$$\begin{aligned} p(x_i|\theta) &= \int \mathcal{N}(x_i|Wz_i + \mu, \Psi) \mathcal{N}(z_i|\mu_0, \Sigma_0) dz_i \\ &= \mathcal{N}(x_i|W\mu_0 + \mu, \Psi + W\Sigma_0 W^T) \end{aligned}$$

- From this, we see that we can set $\mu_0 = 0$ without a loss of generality.
 - we can always absorb $W\mu_0$ into μ .
- Similarly, we can set $\Sigma_0 = I$.
 - we can always emulate a correlated prior by defining a new weight matrix $\tilde{W} = W\Sigma_0^{\frac{1}{2}}$.
- Therefore,

$$\text{cov}[x|\theta] = W\Sigma_0 W^T + \Psi = \left(\tilde{W}\Sigma_0^{-\frac{1}{2}}\right)\Sigma_0\left(\tilde{W}\Sigma_0^{-\frac{1}{2}}\right)^T + \Psi = \tilde{W}\tilde{W}^T + \Psi$$

Factor Analysis

- We thus see that FA approximates the covariance matrix of the visible vector using a **low-rank decomposition**:

$$\mathbf{C} \triangleq \text{cov}[\mathbf{x}] = \mathbf{W}\mathbf{W}^T + \boldsymbol{\Psi}$$

- \mathbf{W} only uses $\mathcal{O}(MD)$ parameters, which allows a flexible compromise between a full covariance Gaussian with $\mathcal{O}(D^2)$ and a diagonal covariance with $\mathcal{O}(D)$ parameters.
- Note that if we did not restrict $\boldsymbol{\Psi}$ to be diagonal, we could trivially set $\boldsymbol{\Psi}$ to a full covariance matrix.

Then we can set $\mathbf{W} = \mathbf{0}$. In this case, latent factors will not be required.



Inference of Latent Factors

- Although FA can be thought of as just a way to define a density on x , it is often used because we hope that the latent factors z will reveal something interesting about the data.
- To do this, we need to compute the posterior over the latent factors.
- We can use Bayes' rule for Gaussians to give

$$\begin{aligned} p(z_i | x_i, \theta) &= \mathcal{N}(z_i | m_i, \Sigma_i) \\ \Sigma_i &\triangleq (\Sigma_0^{-1} + \mathbf{W}\Psi^{-1}\mathbf{W})^{-1} \\ m_i &\triangleq \Sigma_i(\mathbf{W}^T\Psi^{-1}(x_i - \mu) + \Sigma_0^{-1}\mu_0) \end{aligned}$$

- Note that in the FA model, Σ_i is actually independent of i . Computing this matrix takes $\mathcal{O}(M^3 + M^2D)$ time.
- Computing each $m_i = \mathbb{E}[z_i | x_i, \theta]$ takes $\mathcal{O}(M^2 + MD)$ time.
- The m_i are called latent score or latent factors.



Inference of Latent Factors

- Use

$$p(\mathbf{x}_i | \mathbf{z}_i, \boldsymbol{\theta}) = \mathcal{N}(\mathbf{x}_i | \mathbf{W}\mathbf{z}_i + \boldsymbol{\mu}, \boldsymbol{\Psi}),$$

$$p(\mathbf{z}_i) = \mathcal{N}(\mathbf{z}_i | \boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0),$$

Bayes' rule and standard eqs. for linear Gaussian systems to derive the posterior of the latent variables:

$$\begin{aligned} p(\mathbf{z}_i | \mathbf{x}_i, \boldsymbol{\theta}) &= \mathcal{N}(\mathbf{z}_i | \mathbf{m}_i, \boldsymbol{\Sigma}_i), \\ \mathbf{m}_i &= \boldsymbol{\Sigma} (\mathbf{W}^T \boldsymbol{\Psi}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) + \boldsymbol{\Sigma}_0^{-1} \boldsymbol{\mu}_0) \\ \boldsymbol{\Sigma} &= (\boldsymbol{\Sigma}_0^{-1} + \mathbf{W}^{-1} \boldsymbol{\Psi}^{-1} \mathbf{W})^{-1} \end{aligned}$$

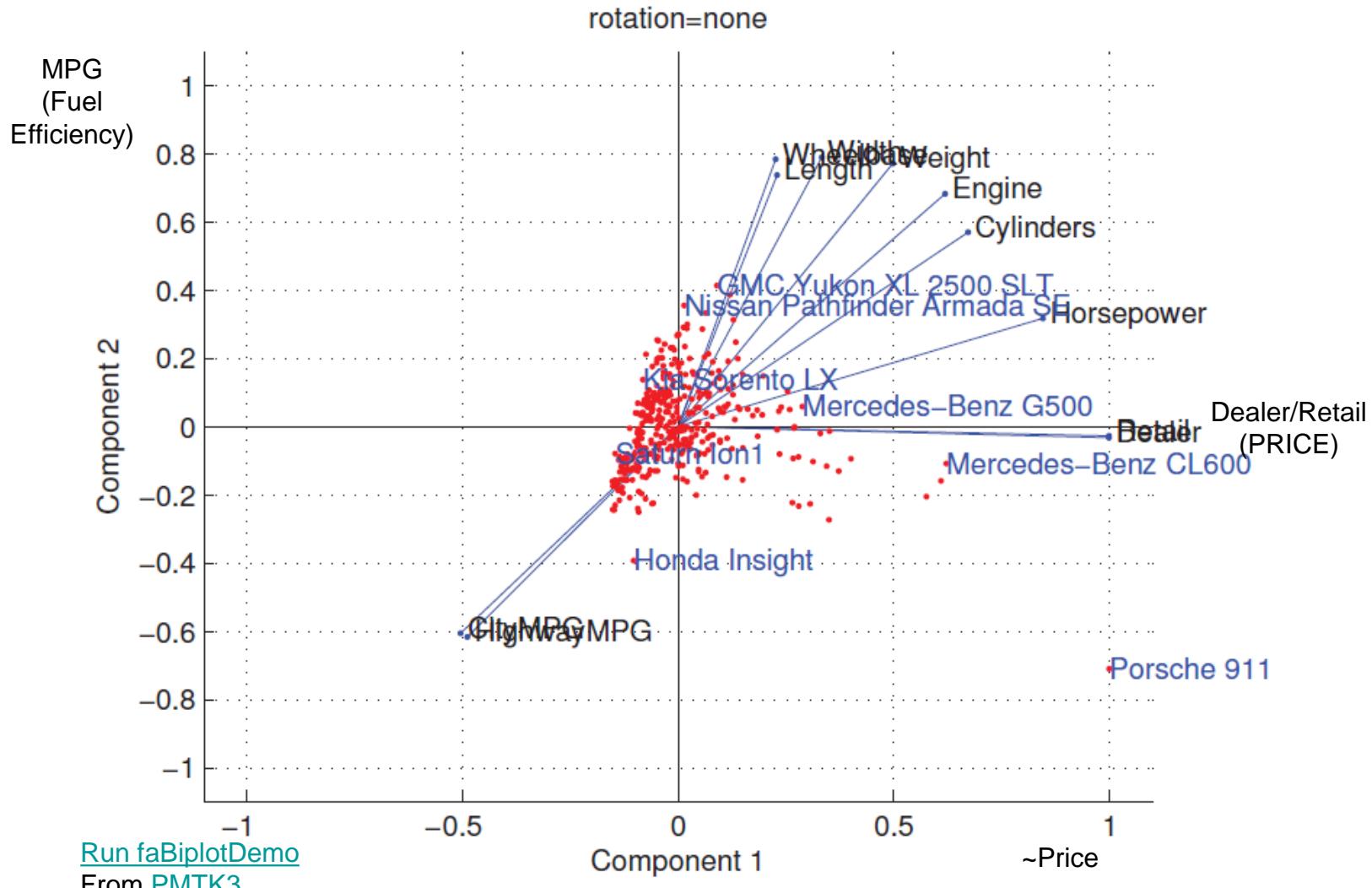
Example Based on Factor Analysis

- Consider a dataset of $D = 11$ variables and $N = 387$ cases describing aspects of cars (engine size, # of cylinders, MPG, etc.).
- We fit a $M = 2$ dimensional model. We plot the $\mathbf{m}_i = \mathbb{E}[\mathbf{z}_i | \mathbf{x}_i, \theta]$ scores as points in \mathbb{R}^2 , to visualize the data.
- We also project unit vectors corresponding to each of the feature dimensions, $\mathbf{e}_1 = (1, 0, \dots, 0)$, $\mathbf{e}_2 = (0, 1, 0, \dots, 0)$, etc. into the low-dimensional space (blue lines).
- We see that the horizontal axis represents price corresponding to the features labeled “dealer” and “retail” (expensive cars on the right).
- The vertical axis represents fuel efficiency (in terms of MPG) versus size: heavy vehicles are less efficient and are higher up, whereas light vehicles are lower down.
- Unfortunately, interpreting latent variable models is fraught with difficulties.



2D Projection of Car Data

- Shalizi, C. (2009). [Cs 36-350 lecture 10: Principal components: mathematics, example, interpretation.](#)



Unidentifiability



Unidentifiability

- Consider an arbitrary rotation R such that: $\tilde{W} = WR$.
- We can immediately see that the likelihood function $p(x_i|\theta) = \mathcal{N}(x_i|\mu, \Psi + WW^T)$ remains the same:

$$\tilde{W} \tilde{W}^T = WRR^TW^T = WW^T$$

- We need to remove $M(M - 1)/2$ DOF since that is the number of orthonormal matrices of size $M \times M$.
- Then the FA model has $D + MD - M(M - 1)/2$ free parameters (excluding the mean), where the first term arises from Ψ . This needs to be $\leq D(D + 1)/2$, i.e. the number of parameters in an unconstrained symmetric covariance matrix.
- This gives us an upper bound on M : $M_{max} = \lfloor D + 0.5(1 -$

Unidentifiability

To address the unidentifiability the following options are available:

□ Making \mathbf{W} orthonormal: This is the case with PCA

- columns arranged in order of decreasing variance

□ Making \mathbf{W} lower triangular: The first feature is then generated only by the 1st latent variable, the second by the first two, etc.

- The number of parameters is equal to M_{max} .
- One needs to properly select the first M visible variables as they affect the latent factors.

- [Lopes, H. and M. West \(2004\). Bayesian model assessment in factor analysis. *Statistica Sinica* 14, 41– 67.](#)



Unidentifiability

- Sparsity promoting priors on the weights: Use sparse factor analysis that forces some entries in W to zero. They may include ℓ_1 regularization, automatic relevance determination or spike and slab priors. But note that this does not ensure a unique MAP estimate.
- Choosing an informative rotation matrix R : Choosing appropriate R to increase interpretable sparse W .
- Using non-Gaussian priors for the latent factors in $p(\mathbf{z}_i)$: This leads to Independent Component Analysis (ICA).

- Zou, H. (2006). [The adaptive Lasso and its oracle properties](#). *J. of the Am. Stat. Assoc.*, 1418–1429.
- Bishop, C. (1999). [Bayesian PCA](#). In *NIPS*.
- Archambeau, C. and F. Bach (2008). [Sparse probabilistic projections](#). In *NIPS*.
- Kaiser, H. (1958). [The varimax criterion for analytic rotation in factor analysis](#), *Psychometrika* 23(3).



Mixtures of Factor Analysers



Mixtures of Factor Analysers

- The FA model assumes that the data lives on a low-dimensional linear manifold.
- In reality, most data is better modeled by some form of low-dimensional curved manifold.
- We can approximate a curved manifold by a piecewise linear manifold
 - This suggests the following model. Let the k -th linear subspace of dimensionality M_k be represented by \mathbf{W}_k , for $k = 1, \dots, K$.
 - Suppose, we have a latent indicator $q_i \in \{1, \dots, K\}$ specifying which subspace we should use to generate the data.
 - We then sample z_i from a Gaussian prior and pass it through the \mathbf{W}_k matrix and add noise:

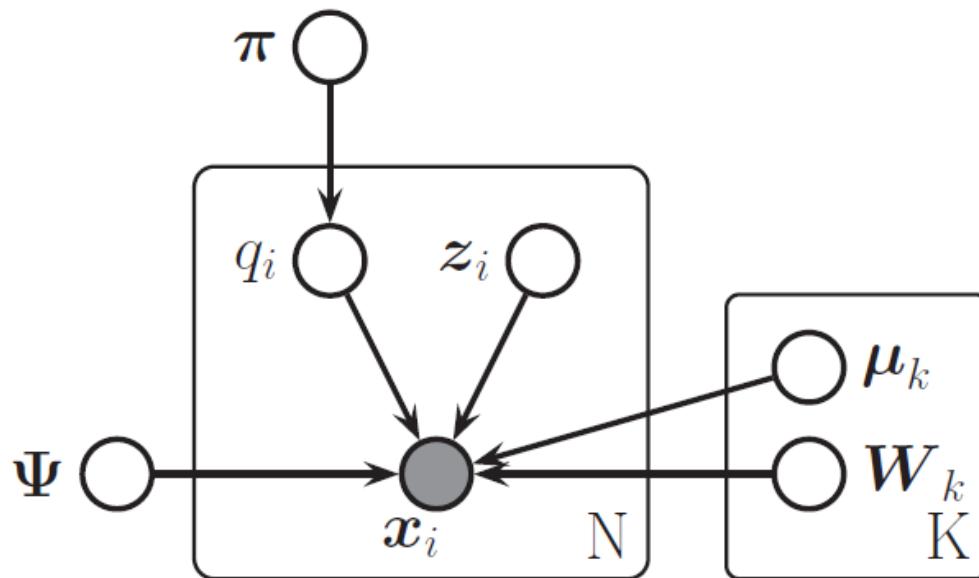
$$\begin{aligned} p(\mathbf{x}_i | \mathbf{z}_i, q_i = k, \boldsymbol{\theta}) &= \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k + \mathbf{W}_k \mathbf{z}_i, \boldsymbol{\Psi}) \\ p(\mathbf{z}_i | \boldsymbol{\theta}) &= \mathcal{N}(z_i | \mathbf{0}, \mathbf{I}) \\ p(q_i | \boldsymbol{\theta}) &= \text{Cat}(q_i | \boldsymbol{\pi}) \end{aligned}$$

- This is called mixture of factor analysers.



Mixtures of Factor Analysers

- A probabilistic graphical model of the mixture of factor analysers is given below.

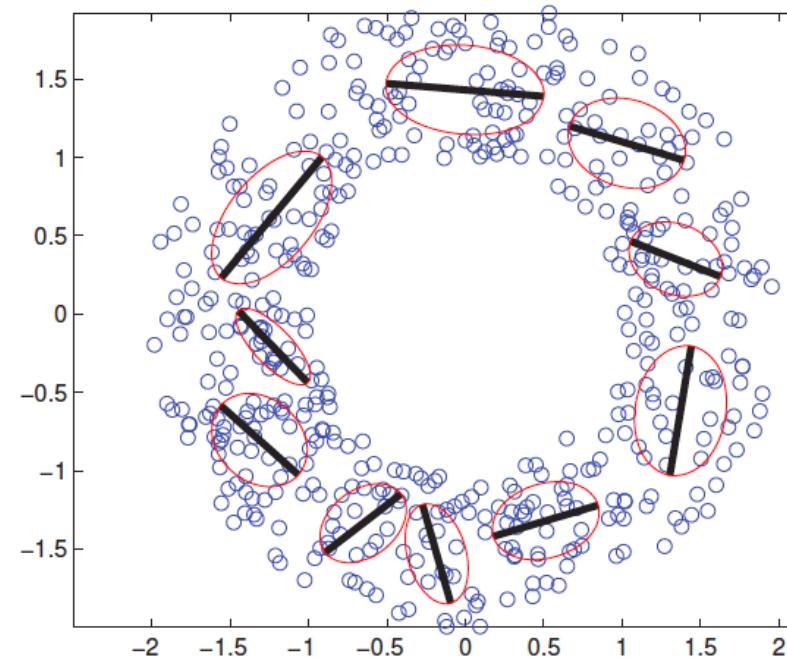
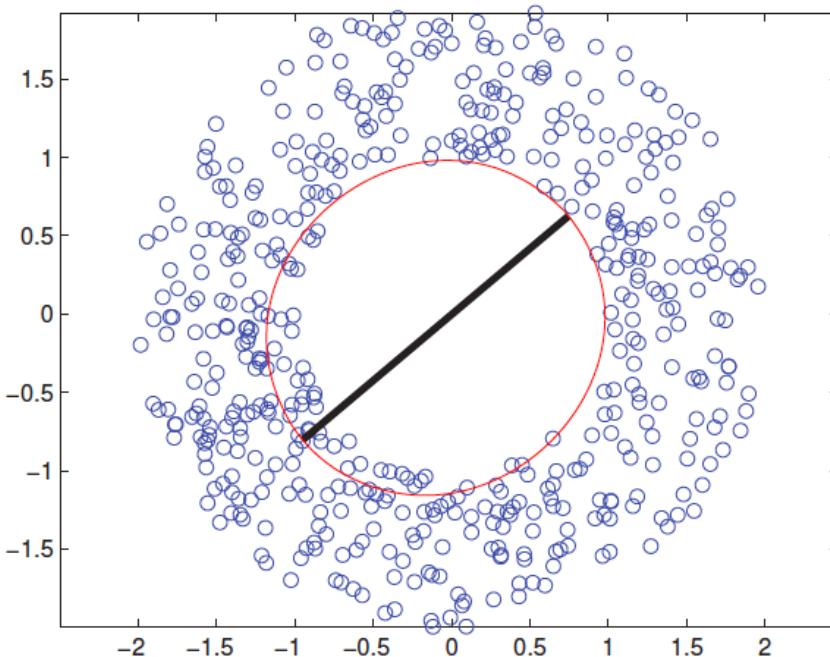


$$\begin{aligned} p(x_i | z_i, q_i = k, \theta) &= \mathcal{N}(x_i | \mu_k + W_k z_i, \Psi) \\ p(z_i | \theta) &= \mathcal{N}(z_i | \mathbf{0}, \mathbf{I}) \\ p(q_i | \theta) &= \text{Cat}(q_i | \pi) \end{aligned}$$

Mixtures of Factor Analysers

- Mixture of 1d PPCAs fit to a dataset, for $K = 1, 10$.

[mixPpcaDemoNetlab](#)
from [PMTK3](#)



- Think of this is as a low-rank version of a mixture of Gaussians
- The model needs $\mathcal{O}(KMD)$ parameters instead of the $\mathcal{O}(KD^2)$ parameters needed for a mixture of full-covariance Gaussians.
- This can reduce overfitting and thus **MFA** is a good generic density model for high-dimensional real-valued data.

Latent Models

- PCA arises as the Maximum Likelihood solution to a linear Gaussian Latent Variable Model.
- Probabilistic formulation allows:
 - Use of EM for parameter estimation
 - Using mixtures of PCA models
 - A Bayesian formalism in addition allows computing the number of principal components directly from the data.
- Non Gaussian latent variable models lead to Independent Component Analysis.
- Models can also be considered with a non-linear relation between latent and observed variables.

Principal Component Analysis



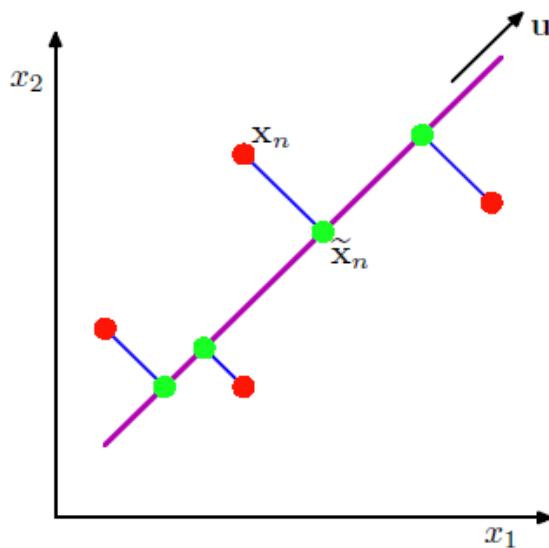
Principal Component Analysis

- Principal component analysis (PCA) – also known as Karhunen-Loeve expansion - is used widely for (Jolliffe, 2002)
 - Dimensionality reduction
 - Lossy data compression
 - Feature selection and
 - Data Visualization
 - Principal component analysis seeks a lower dimensionality linear space (principal subspace) such that
 - The orthogonal projection of the data points onto this subspace **maximizes the variance of the projected points** (Hotelling, 1933)
 - An alternative definition of PCA is based on **minimizing the sum-of-squares of the projection errors** (Pearson, 1901)
-
- Jolliffe, I. T. (2002). Principal Component Analysis (Second ed.). Springer.
 - Hotelling, H. (1933). Analysis of a complex of statistical variables into principal components. *Journal of Educational Psychology* **24**, 417–441.
 - Pearson, K. (1901). On lines and planes of closest fit to systems of points in space. *The London, Edinburgh and Dublin Philosophical Magazine and Journal of Science, Sixth Series* **2**, 559–572.



Principal Component Analysis

- PCA seeks a space of lower dimensionality (magenta line) such that:
 - (1) the orthogonal projection of the data points (red dots) onto this subspace **maximizes the variance of the projected points** (green dots).
 - (2) minimizing the sum-of-squares of the projection errors (blue lines)



- J. Shlens (2005). [A Tutorial on Principal Component Analysis](#)

PCA: Maximum Variance Formulation

- Consider a data set of observations $\{x_n\}, n = 1, \dots, N$, where x_n has dimensionality D .
- Our goal is to project the data onto a space having dimensionality $M < D$ (principal subspace) while maximizing the variance of the projected data (M here is fixed).
- Let $\{u_i\}, i = 1, \dots, M$ the basis vectors ($(D \times 1)$ vectors) of the principal subspace.
- We define the sample mean ($(D \times 1)$ vector) by:

$$\bar{x} = \frac{1}{N} \sum_{n=1}^N x_n$$

and the sample covariance matrix ($(D \times D)$ matrix) by:

$$S = \frac{1}{N} \sum_{n=1}^N (x_n - \bar{x})(x_n - \bar{x})^T$$



PCA with 1D Principal Subspace

- Consider first the projection onto 1D space ($M = 1$).
- We define the direction of this subspace using vector \mathbf{u}_1 , which we choose to be a unit vector: $\mathbf{u}_1^T \mathbf{u}_1 = 1$
- Each data point \mathbf{x}_n is then projected onto the scalar $\mathbf{u}_1^T \mathbf{x}_n$
- The mean of the projected data is given as: $\mathbf{u}_1^T \bar{\mathbf{x}}$
- Similarly, the variance of the projected data is:

$$\frac{1}{N} \sum_{n=1}^N \left\{ \mathbf{u}_1^T \mathbf{x}_n - \mathbf{u}_1^T \bar{\mathbf{x}} \right\}^2 = \frac{1}{N} \sum_{n=1}^N \left\{ \mathbf{u}_1^T \mathbf{x}_n - \mathbf{u}_1^T \bar{\mathbf{x}} \right\} \left\{ \mathbf{u}_1^T \mathbf{x}_n - \mathbf{u}_1^T \bar{\mathbf{x}} \right\}^T = \mathbf{u}_1^T \mathbf{S} \mathbf{u}_1$$

- Key idea of PCA: Maximize the projected variance $\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1$ with respect to \mathbf{u}_1 under the constraint $\mathbf{u}_1^T \mathbf{u}_1 = 1$

PCA with 1D Principal Subspace

- Introduce the Lagrange multiplier λ_1 and define the unconstrained maximization of

$$\max_{\mathbf{u}_1, \lambda_1} \left\{ \mathbf{u}_1^T \mathbf{S} \mathbf{u}_1 + \lambda_1 (1 - \mathbf{u}_1^T \mathbf{u}_1) \right\}$$

- We can see immediately that the solution satisfies:

$$\mathbf{S} \mathbf{u}_1 = \lambda_1 \mathbf{u}_1$$

- \mathbf{u}_1 must be an eigenvector of \mathbf{S} with eigenvalue λ_1
- From the eigen-problem note that the variance of the projected data is $\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1 = \lambda_1$ so λ_1 needs to be the largest eigenvalue of \mathbf{S} .
- \mathbf{u}_1 is called the first principal component.

PCA with 1D Principal Subspace

- Additional Principal Components: maximize the projected variance amongst all possible directions orthogonal to those already considered.
- Using induction, you can show: For an M -dimensional projection space,

The optimal linear projection for which the variance of the projected data is maximized is defined by the M eigenvectors $\mathbf{u}_1, \dots, \mathbf{u}_M$ of the data covariance matrix \mathbf{S} corresponding to the largest eigenvalues $\lambda_1, \dots, \lambda_M$.
- The computational cost of finding the M principal components (i.e. finding the first M eigenvalues and eigenvectors of \mathbf{S}) is $\mathcal{O}(MD^2)$.



Total Variance

- PCA becomes the calculation of the eigenvectors of the data covariance matrix corresponding to the largest M eigenvalues.
- $\sum_{i=1}^D \lambda_i$ is the **total variance**.
- The contribution to the total variance by one component u_i is measured as:

$$\frac{\lambda_i}{\sum_{i=1}^D \lambda_i}$$



Proof by Induction: M -Dimensional Subspace

- We have shown our result for $M = 1$. Assume it is valid for projection space of dimensionality M . We will show that it is then valid for projection space of dimensionality $M + 1$.
- We need to maximize the following:

$$\max_{\mathbf{u}_{M+1}} \left\{ \mathbf{u}_{M+1}^T S \mathbf{u}_{M+1} + \lambda_{M+1} \left(1 - \mathbf{u}_{M+1}^T \mathbf{u}_{M+1} \right) + \sum_{i=1}^M \eta_i \mathbf{u}_{M+1}^T \mathbf{u}_i \right\}$$

- The stationary points satisfy:

$$0 = 2S\mathbf{u}_{M+1} - 2\lambda_{M+1}\mathbf{u}_{M+1} + \sum_{i=1}^M \eta_i \mathbf{u}_i$$

- Taking dot product with $\mathbf{u}_j^T, j = 1, \dots, M$ gives $\eta_j = 0$. Thus $S\mathbf{u}_{M+1} = \lambda_{M+1}\mathbf{u}_{M+1} \Rightarrow \mathbf{u}_{M+1}$ eigenvector of S with the variance

in this direction $\mathbf{u}_{M+1}^T S \mathbf{u}_{M+1} = \lambda_{M+1}$

- λ_{M+1} needs to be the next largest eigenvalue after the first M ones \rightarrow The argument by induction holds for any $M \leq D$.

PCA: Minimum Error Formulation

- This is based on projection error minimization.
- Consider a D -dimensional basis $\{\mathbf{u}_i\}$ where $i = 1, \dots, D$ satisfying

$$\mathbf{u}_1^T \mathbf{u}_j = \delta_{ij}$$

- Each data point \mathbf{x}_n can be represented by:

$$\mathbf{x}_n = \sum_{i=1}^D a_{ni} \mathbf{u}_i, \text{ where: } a_{ni} = \mathbf{x}_n^T \mathbf{u}_i$$

- We approximate \mathbf{x}_n using a representation with $M < D$ variables z_{ni} via a projection onto a lower-dim subspace:

$$\tilde{\mathbf{x}}_n = \sum_{i=1}^M z_{ni} \mathbf{u}_i + \sum_{i=M+1}^D b_i \mathbf{u}_i$$

where the $\{z_{ni}\}$ depend on the particular data point, whereas the $\{b_i\}$ are the same constants for all data points.

- We are free to choose the $\{\mathbf{u}_i\}$, $\{z_{ni}\}$, and $\{b_i\}$

PCA: Minimum Error Formulation

- Key idea of PCA: Minimize the distortion J introduced by the dimensionality reduction:

$$\min_{z_{nj}, b_j, \mathbf{u}_i} J = \frac{1}{N} \sum_{i=1}^N \| \mathbf{x}_n - \tilde{\mathbf{x}}_n \|^2, \quad \tilde{\mathbf{x}}_n = \sum_{i=1}^M z_{ni} \mathbf{u}_i + \sum_{i=M+1}^D b_i \mathbf{u}_i$$

- Setting the derivatives wrt $\{z_{ni}\}$, and $\{b_i\}$ to zero and using $\mathbf{u}_1^T \mathbf{u}_j = \delta_{ij}$ gives:

$$z_{nj} = \mathbf{x}_n^T \mathbf{u}_j, \quad j = 1, \dots, M$$

$$b_j = \bar{\mathbf{x}}^T \mathbf{u}_j, \quad j = M+1, \dots, D, \text{ where } \bar{\mathbf{x}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n$$

- If we now substitute these expressions for $\{z_{ni}\}$, and $\{b_i\}$

$$\mathbf{x}_n - \tilde{\mathbf{x}}_n = \sum_{i=1}^D (\mathbf{x}_n^T \mathbf{u}_i) \mathbf{u}_i - \sum_{i=1}^M (\mathbf{x}_n^T \mathbf{u}_i) \mathbf{u}_i - \sum_{i=M+1}^D (\bar{\mathbf{x}}^T \mathbf{u}_i) \mathbf{u}_i = \sum_{i=M+1}^D \{(\mathbf{x}_n - \bar{\mathbf{x}})^T \mathbf{u}_i\} \mathbf{u}_i$$

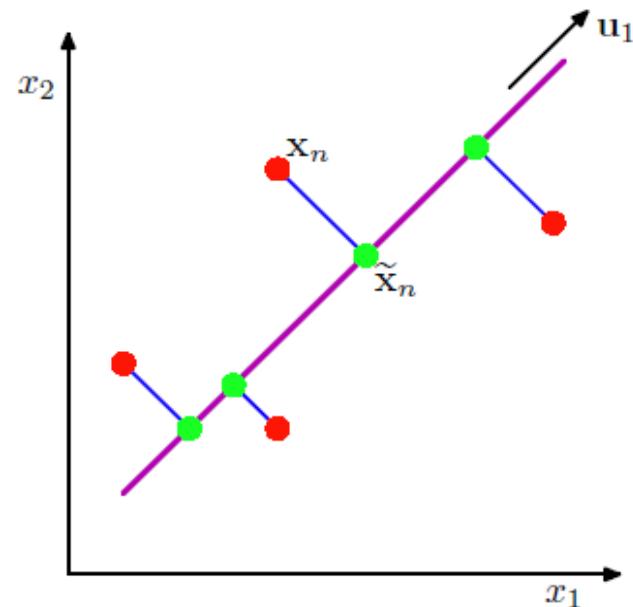
PCA: Minimum Error Formulation

$$x_n - \tilde{x}_n = \sum_{i=M+1}^D \{(x_n - \bar{x})^T u_i\} u_i$$

- We see that the displacement vector $x_n - \tilde{x}_n$ lies in the space orthogonal to the principal subspace, i.e. is a linear combination of $\{u_i\}$ for $i = M + 1, \dots, D$.

- We can also see this geometrically:

- The projected points \tilde{x}_n must lie within the principal subspace, but can move them freely within that subspace.
- The minimum error is then obtained by the orthogonal projection.



PCA: Minimum Error Formulation

- We obtain an expression for the distortion measure J as a function purely of the $\{\mathbf{u}_i\}$ in the form

$$\mathbf{x}_n - \tilde{\mathbf{x}}_n = \sum_{i=M+1}^D \left\{ (\mathbf{x}_n - \bar{\mathbf{x}})^T \mathbf{u}_i \right\} \mathbf{u}_i$$
$$J = \frac{1}{N} \sum_{n=1}^N \sum_{i=M+1}^D \left(\mathbf{x}_n^T \mathbf{u}_i - \bar{\mathbf{x}}^T \mathbf{u}_i \right)^2 = \sum_{i=M+1}^D \mathbf{u}_i^T \mathbf{S} \mathbf{u}_i$$
$$J = \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}_n - \tilde{\mathbf{x}}_n\|^2$$

- The minimum of J is obtained when $\{\mathbf{u}_i\}, i = M + 1, \dots, D$ are the eigenvectors of \mathbf{S} associated to the smallest eigenvalues.
- Consider $D = 2$ and $M = 1$. Let $\lambda_1 > \lambda_2$. The principal subspace is aligned with the eigenvector with the larger eigenvalue λ_1 , and the min value of $J = \lambda_2$ is obtained by choosing \mathbf{u}_2 corresponding to λ_2 .

PCA: Minimum Error Formulation

- The distortion measure is given as

$$J = \sum_{i=M+1}^D \lambda_i$$

- Obtain the minimum value of J by selecting these eigenvectors corresponding to the $D - M$ smallest eigenvalues.
- The eigenvectors defining the principal subspace are those corresponding to the M largest eigenvalues.
- The two approaches discussed
 - Maximum Variance and
 - Minimum Errorboth lead to the same main result of the PCA.



Compression of the Original Data Set

- Using the earlier equations, we can derive:

$$\begin{aligned}\tilde{\mathbf{x}}_n &= \sum_{i=1}^M z_{ni} \mathbf{u}_i + \sum_{i=M+1}^D b_i \mathbf{u}_i \\ z_{nj} &= \mathbf{x}_n^T \mathbf{u}_j, j = 1, \dots, M & \Rightarrow \tilde{\mathbf{x}}_n &= \sum_{i=1}^M (\mathbf{x}_n^T \mathbf{u}_i) \mathbf{u}_i + \sum_{i=M+1}^D (\bar{\mathbf{x}}^T \mathbf{u}_i) \mathbf{u}_i \Rightarrow \\ b_j &= \bar{\mathbf{x}}^T \mathbf{u}_j, j = M+1, \dots, D\end{aligned}$$

Data Compression: $\tilde{\mathbf{x}}_n = \bar{\mathbf{x}} + \sum_{i=1}^M (\mathbf{x}_n^T \mathbf{u}_i - \bar{\mathbf{x}}^T \mathbf{u}_i) \mathbf{u}_i$

- In deriving the last equation, we used the completeness of \mathbf{u}_i , i.e.

$$\bar{\mathbf{x}} = \sum_{i=1}^D (\bar{\mathbf{x}}^T \mathbf{u}_i) \mathbf{u}_i$$

- This clearly shows the compression of the data set: From the D -dimensional \mathbf{x}_n to the M -dimensional vector with components $(\mathbf{x}_n^T \mathbf{u}_i - \bar{\mathbf{x}}^T \mathbf{u}_i)$.

PCA Reconstruction

- To reconstruct the data in the original D -dimensional space from a representation in the M -dimensional principal subspace, we simply use:

$$\tilde{\mathbf{x}}_n = \bar{\mathbf{x}} + \sum_{i=1}^M (\mathbf{x}_n^T \mathbf{u}_i - \bar{\mathbf{x}}^T \mathbf{u}_i) \mathbf{u}_i$$

- If $M = D$, there is no dimensionality reduction but a rotation to align with the principal components e.g from $\{\mathbf{x}_{n_1}, \dots, \mathbf{x}_{n_D}\}$ to $\{\mathbf{a}_{n_1}, \dots, \mathbf{a}_{n_D}\}$ where:

$$\mathbf{x}_n = \sum_{i=1}^{M=D} a_{ni} \mathbf{u}_i$$



PCA and Singular Value Decomposition

- We have defined the solution to PCA in terms of eigenvectors of the empirical covariance matrix (here denoted as $\widehat{\Sigma}$). The same solution can be obtained based on SVD on the **standardized data matrix X** .

U, V : orthonormal

$\sigma_1, \sigma_2, \dots, \sigma_D$: singular values (≥ 0)

$$\begin{matrix} & D \\ N & \end{matrix} = \begin{matrix} D & N - D \\ \boxed{} & \boxed{} \\ \boxed{} & \boxed{} \end{matrix} \quad \begin{matrix} D \\ \sigma_1 & \dots & \sigma_D \\ \boxed{0} & \boxed{0} & \boxed{0} \end{matrix} \quad \begin{matrix} D \\ \boxed{} & \boxed{} & \boxed{} \end{matrix} \quad \begin{matrix} D \\ \boxed{} & \boxed{} & \boxed{} \end{matrix}$$
$$X = U S V^T$$

- The shaded entries in S and the off diagonal terms are zero. So we don't need to compute the shaded entries of U and S .
- It can be shown that **the eigenvectors of $X^T X$ are equal to V** (right singular vectors of X) and the eigenvalues of $X^T X$ equal to $D = S^2$ (squared singular values) $X^T X = V S U^T U S V^T = V S^2 V^T \rightarrow (X^T X)V = VD$
- Similarly: $XX^T = USV^T VSU^T = US^2 U^T \rightarrow (XX^T)U = UD$

PCA and Singular Value Decomposition

- If the singular values die off quickly, we can produce a rank M approximation of \mathbf{X} as follows:

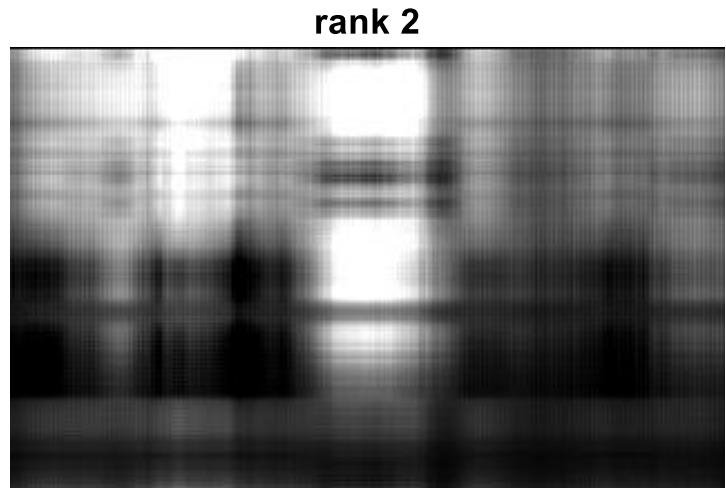
$$\begin{array}{c} D \\ \boxed{\phantom{\text{D}}} \\ N \end{array} \simeq \begin{array}{c} M \\ \boxed{\phantom{\text{M}}} \\ \sigma_1 \\ \vdots \\ \sigma_M \end{array} \begin{array}{c} D \\ \boxed{\phantom{\text{D}}} \\ M \end{array}$$
$$\mathbf{X} \approx \mathbf{U}_M \mathbf{S}_M \mathbf{V}_M^T$$

- The number of parameters for this truncated SVD are: $NM + M + MD$



PCA and Singular Value Decomposition

- Consider the image shown of size 200×320 (rank 200) and the SVD approximations of rank 2, 5, and 20.

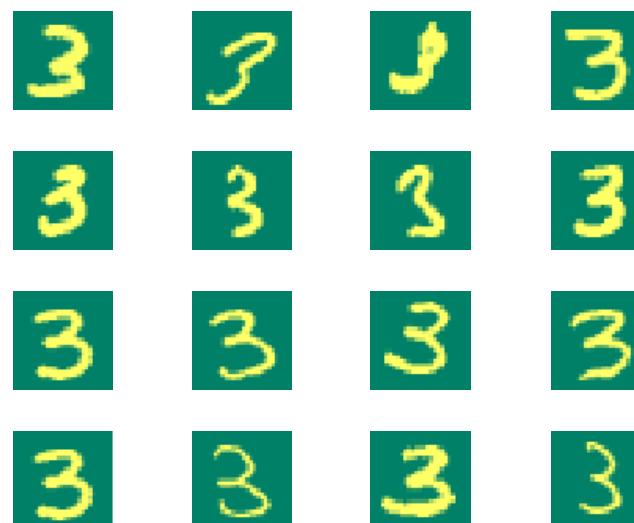


Run [svdImageDemo](#)
From [PMTK3](#)



Off-line Digit Images: An Example of PCA

- Let us apply PCA in the digit images discussed earlier.
- Consider a synthetic data set obtained by taking one of the digit images and creating multiple copies ($N = 10,000$) in each of which the digit has undergone a random displacement and rotation.
- The resulting images each have $D = 28 \times 28 = 784$ pixels.



MatLab Code



Off-line Digit Images: PCA Algorithm

- Our dataset is $X \in D \times N$, where $D = 784$ (pixels) and $N = 10,000$ (images).
- Each column of the matrix X is a data point (image) x_n , ($n = 1, 2, \dots, 10,000$) of size 784.
- We define the centered covariance matrix as:

$$S = \frac{1}{N} \sum_{n=1}^N (x_n - \bar{x})(x_n - \bar{x})^T$$

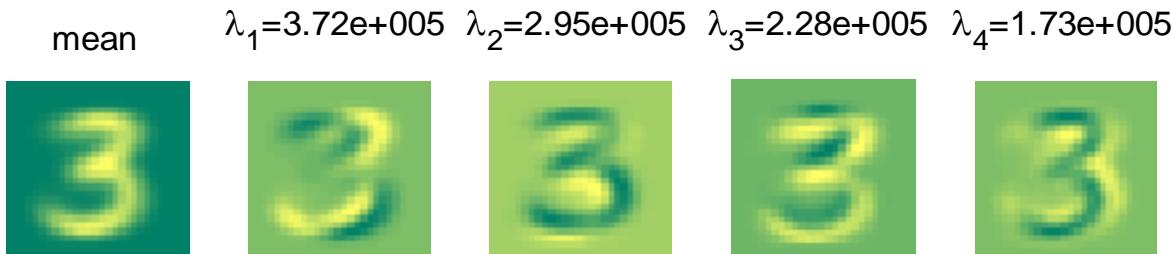
- To compute the principal components, we solve the following eigenvalue problem:

$$Su = \lambda u$$



Off-line Digit Images: PCA Results

- Using the algorithm stated before, we apply PCA on the off-line digits dataset.
- Since each eigenvector of S is a vector in the original D -dimensional space, we can represent the eigenvectors as images of the same size as the data points (see Fig. below)
- Below is the mean vector \bar{x} along with the first four PCA eigenvectors $\mathbf{u}_1, \dots, \mathbf{u}_4$ for the digit 3 from the off-line digits data set, together with the corresponding eigenvalues.



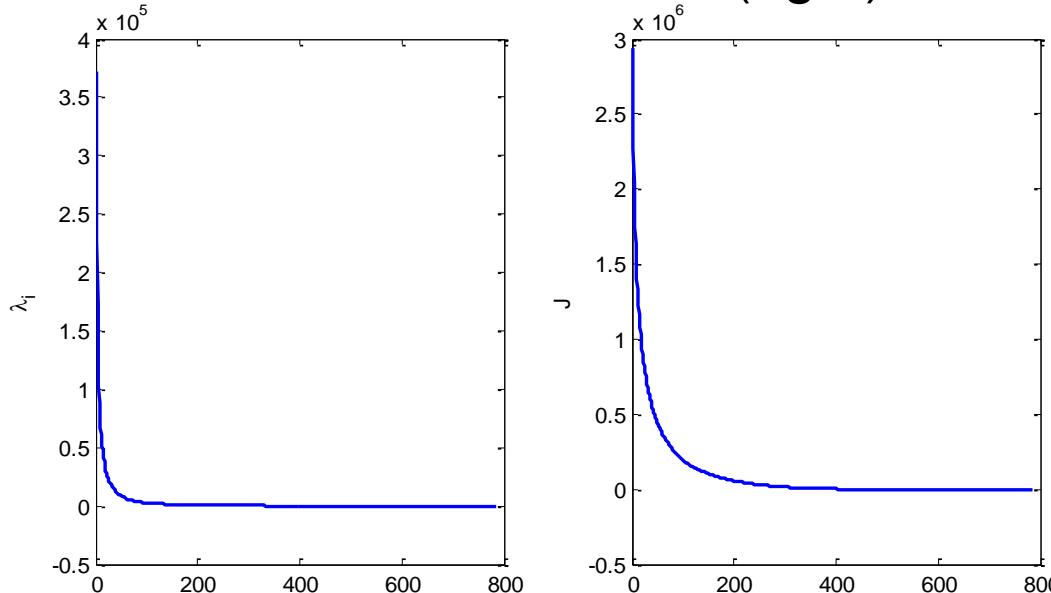
[MatLab Code](#)

Off-line Digit Images: PCA Results

- The spectrum of eigenvalues, sorted into decreasing order, is shown below (left).
- The distortion measure J associated with choosing a particular value of M is given by

$$J = \sum_{i=M+1}^D \lambda_i$$

and is plotted for different values of M (right).



[MatLab Code](#)

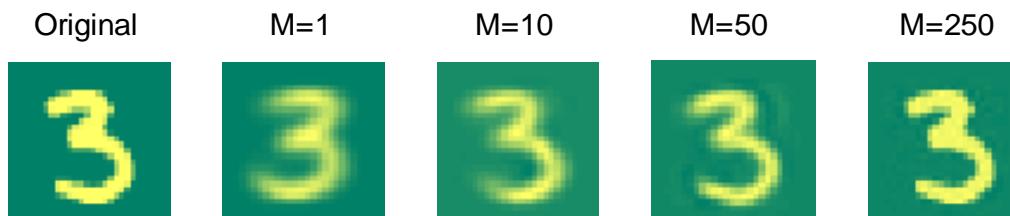


Off-line Digit Images: PCA Reconstruction

- Below shows an original example of digit 3 from the off-line digits data set together with its PCA reconstructions obtained by retaining M principal components for various values of M .

$$\tilde{x}_n = \bar{x} + \sum_{i=1}^M (x_n^T u_i - \bar{x}^T u_i) u_i$$

- As M increases the reconstruction becomes more accurate and would become perfect when $M = D = 28 \times 28 = 784$.



[MatLab Code](#)



Off-line Digit Images: PCA Reconstruction



mean

principal basis 1



Using 5 bases



Using 10 bases



principal basis 2



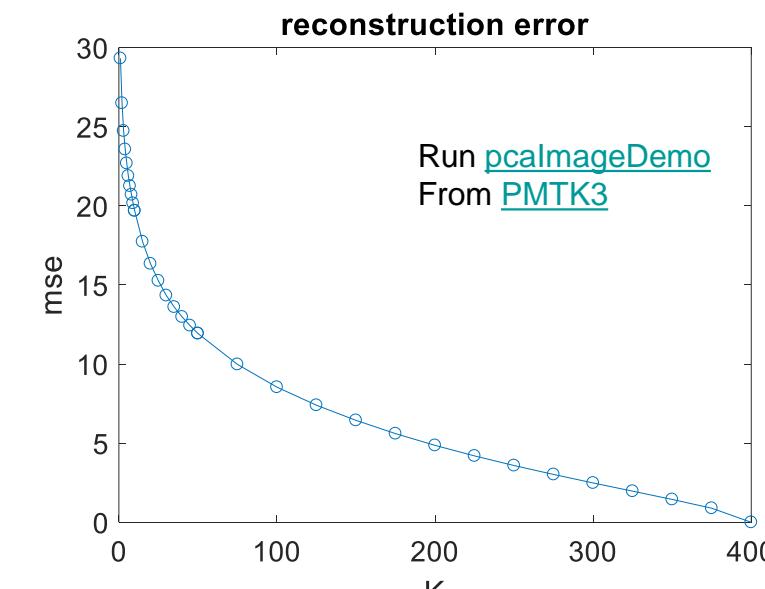
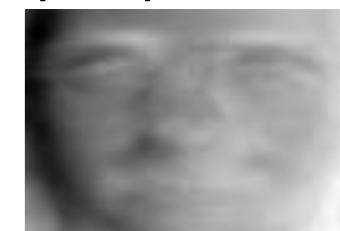
principal basis 3



Using 20 bases

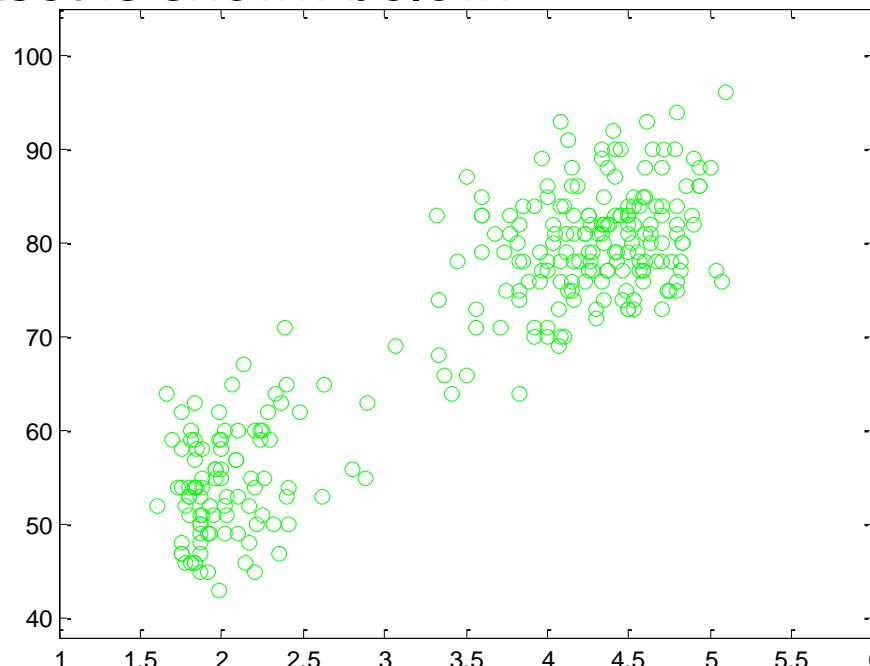


Using 400 bases



Old Faithful Data Set: Whitening

- ❑ Another application of PCA is data pre-processing (standardizing certain of its properties). This is applied here to the old faithful dataset.
- ❑ Here the goal is the transformation of a data set in order to standardize the data (zero mean/unit variance).
- ❑ The original dataset is shown below:



MatLab Code



Old Faithful Data Set: Standardizing

- Standardizing the data is useful for successful application of pattern recognition techniques.
- It is essential when the original variables are measured in various different units or have different scaling.
- In the Old Faithful data set, the time between eruptions is typically an order of magnitude greater than the duration of an eruption.
- Before applying classification techniques (e.g. K –Means) we perform linear re-scaling of the **individual variables** such that **each variable has zero mean and unit variance**.
- This is known as **standardizing the data**.

Old Faithful Data Set: Standardizing

- The covariance matrix for the standardized data has components

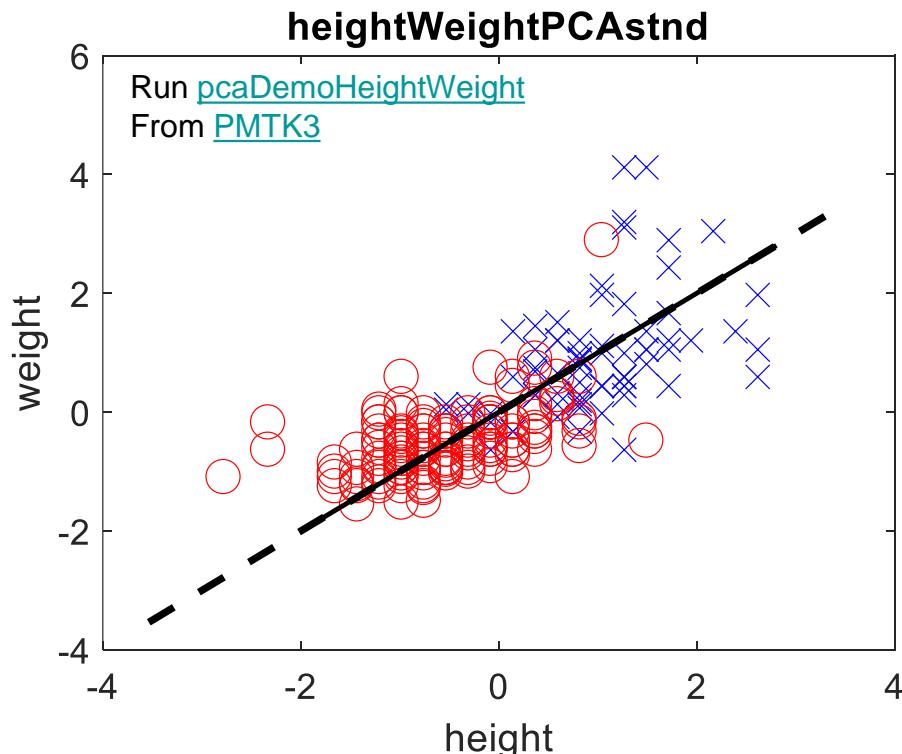
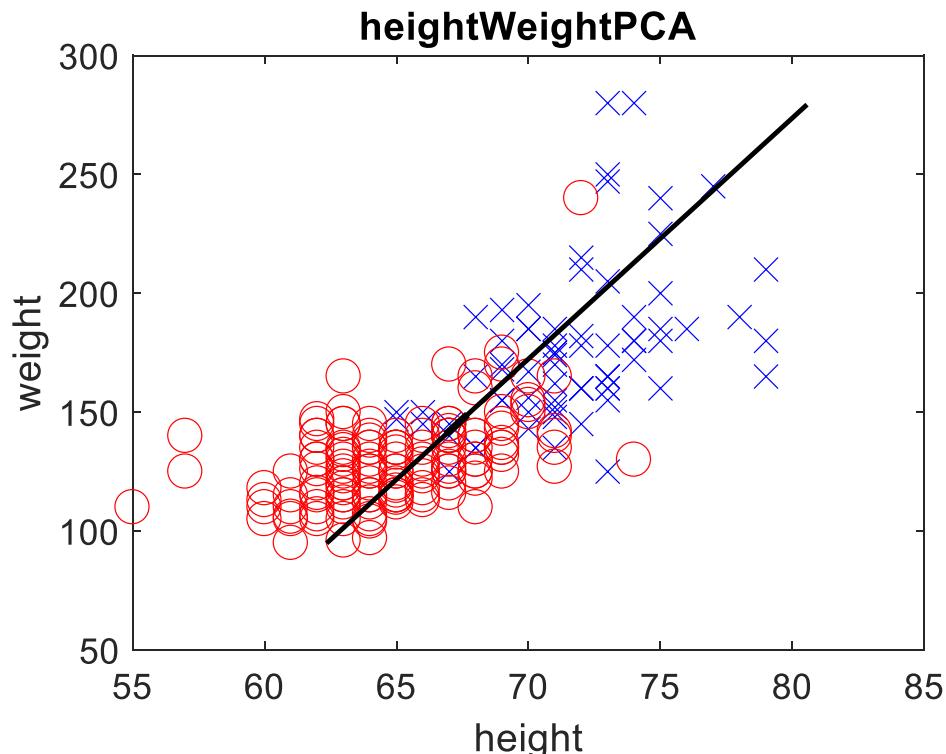
$$\rho_{ij} = \frac{1}{N} \sum_{n=1}^N \frac{x_{ni} - \bar{x}_i}{\sigma_i} \frac{x_{nj} - \bar{x}_j}{\sigma_j}, \sigma_i = \text{std of } x_i$$

- This is the correlation matrix of the original data.
- Two components x_i and x_j of the data are
 - ✓ perfectly correlated, then $\rho_{ij} = 1$, and
 - ✓ uncorrelated, then $\rho_{ij} = 0$.



Old Faithful Data Set: PCA Example

- It should be standard practice to standardize the data first before performing PCA. This is equivalent to working with correlation matrices instead of covariance matrices.
- (Left) PCA for raw data (Right) PCA of standardized data (individual variables having zero mean and unit variance).



Old Faithful Data Set: Whitening

- With PCA we can normalize the data to give them zero mean and **unit covariance (different variables become decorrelated)**.
- Consider the key eigenvalue problem in PCA in a matrix form:

$$\mathbf{S}\mathbf{U} = \mathbf{U}\mathbf{L}, \quad \mathbf{L} = \text{diag}(\lambda_1, \dots, \lambda_D), \quad \mathbf{U} = [\mathbf{u}_1 \dots \mathbf{u}_D] \text{ (orthogonal)}$$

- For each data point \mathbf{x}_n , define a transformed value as:

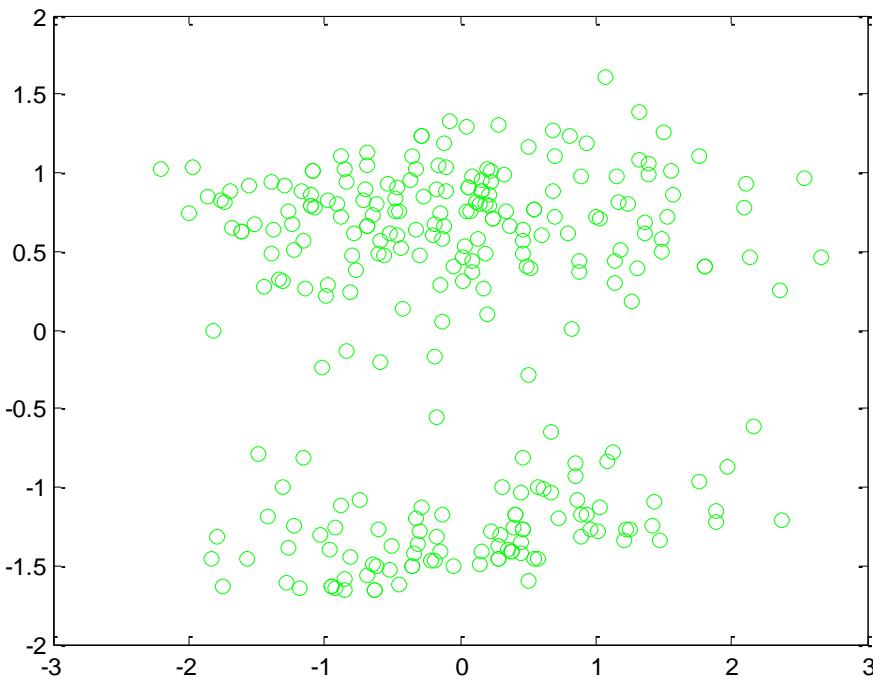
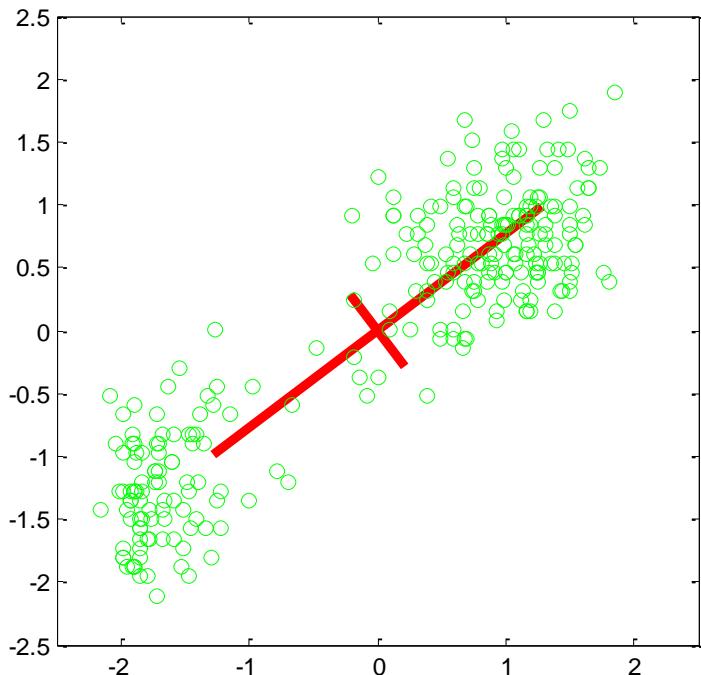
Whitening of the data : $\mathbf{y}_n = \mathbf{L}^{-1/2} \mathbf{U}^T (\mathbf{x}_n - \bar{\mathbf{x}})$

- The set $\{\mathbf{y}_n\}$ has zero mean and its covariance is the identity:

$$\frac{1}{N} \sum_{n=1}^N \mathbf{y}_n \mathbf{y}_n^T = \frac{1}{N} \sum_{n=1}^N \mathbf{L}^{-1/2} \mathbf{U}^T (\mathbf{x}_n - \bar{\mathbf{x}}) (\mathbf{x}_n - \bar{\mathbf{x}})^T \mathbf{U} \mathbf{L}^{-1/2} = \mathbf{L}^{-1/2} \mathbf{U}^T \mathbf{S} \mathbf{U} \mathbf{L}^{-1/2} = \mathbf{I}$$

Old Faithful Data Set: PCA Example

- Left: Standardizing individual variables to zero mean and unit variance. The principal axes of the normalized set are shown for the range $\pm\lambda_i^{1/2}$ (variables still correlated)
- Right: Whitening of the data (zero mean, unit covariance)

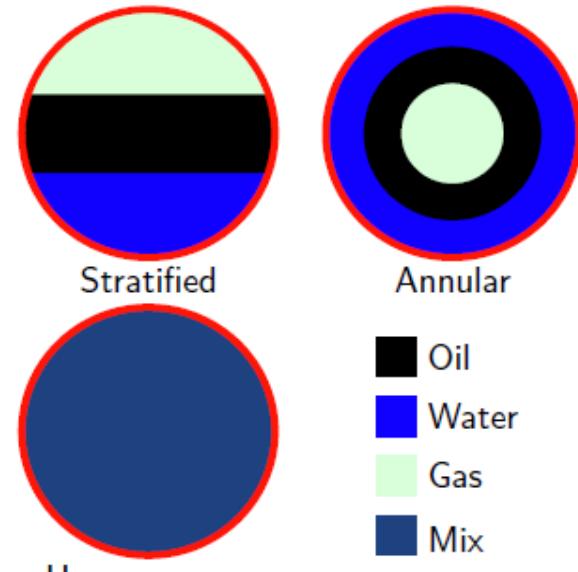


MatLab Code



Oil Flow Data Set

- Consider the oil flow data set.
- For a given geometrical configuration of the gas, water, and oil phases) there are only 2 DOF of variability
 - ✓ the fraction of oil in the pipe and
 - ✓ the fraction of water (the fraction of gas then being determined).
- The data space comprises 12 measurements.
- However, the points in the data set lie close to a 2D manifold embedded within this 12th dimensional space.

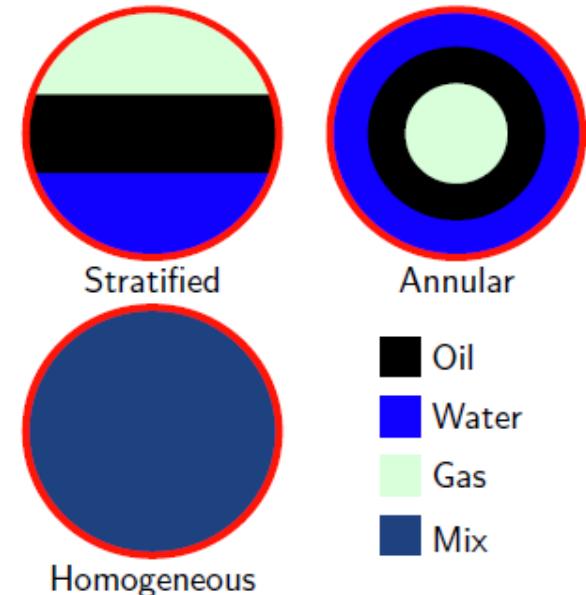


Flow Configurations

Red : 'laminar'
Blue: 'homogeneous'
Green: 'annular'

Oil Flow Data Set

- The manifold comprises **several distinct segments** corresponding to the different flow regimes, **each such segment being a (noisy) continuous 2D manifold**.
- For data compression, there is merit in exploiting this manifold structure.

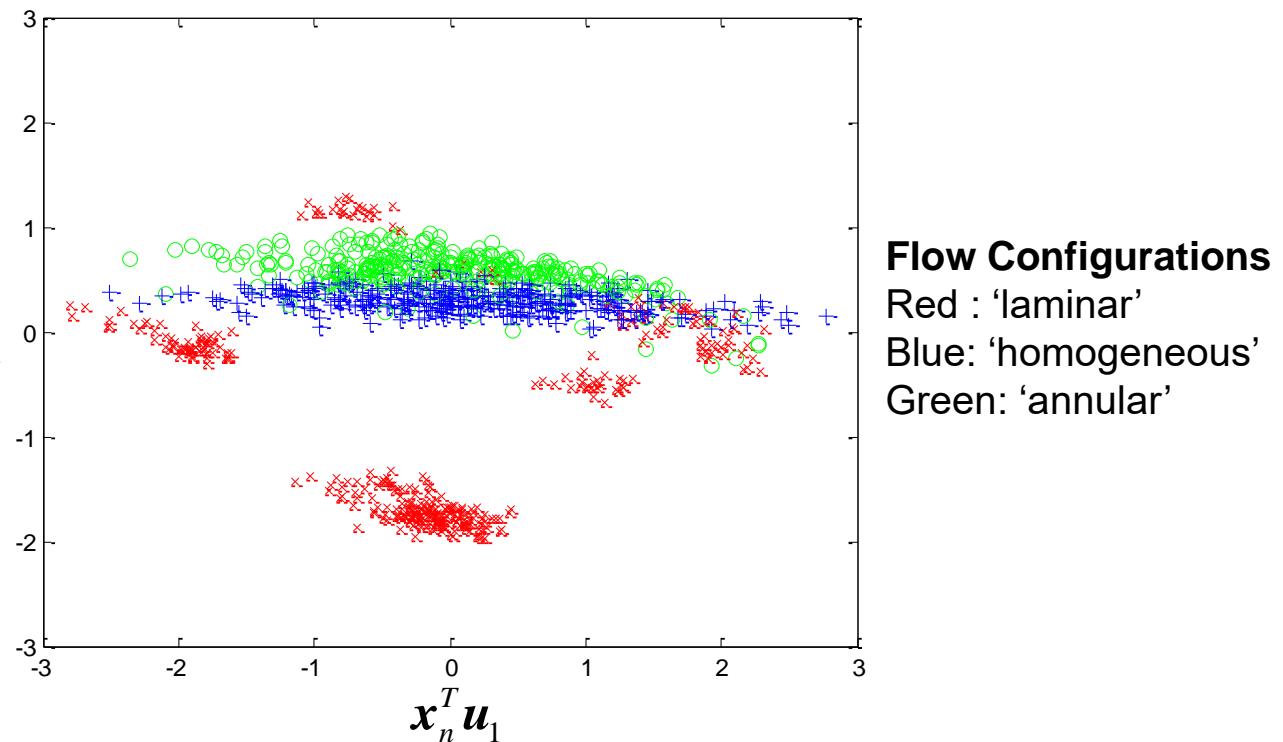
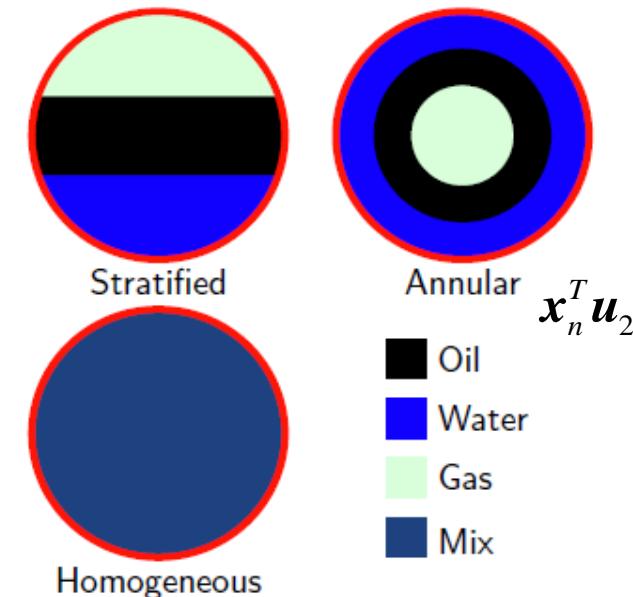


Flow Configurations

Red : 'laminar'
Blue: 'homogeneous'
Green: 'annular'

Oil Flow Data Set: PCA for Visualization

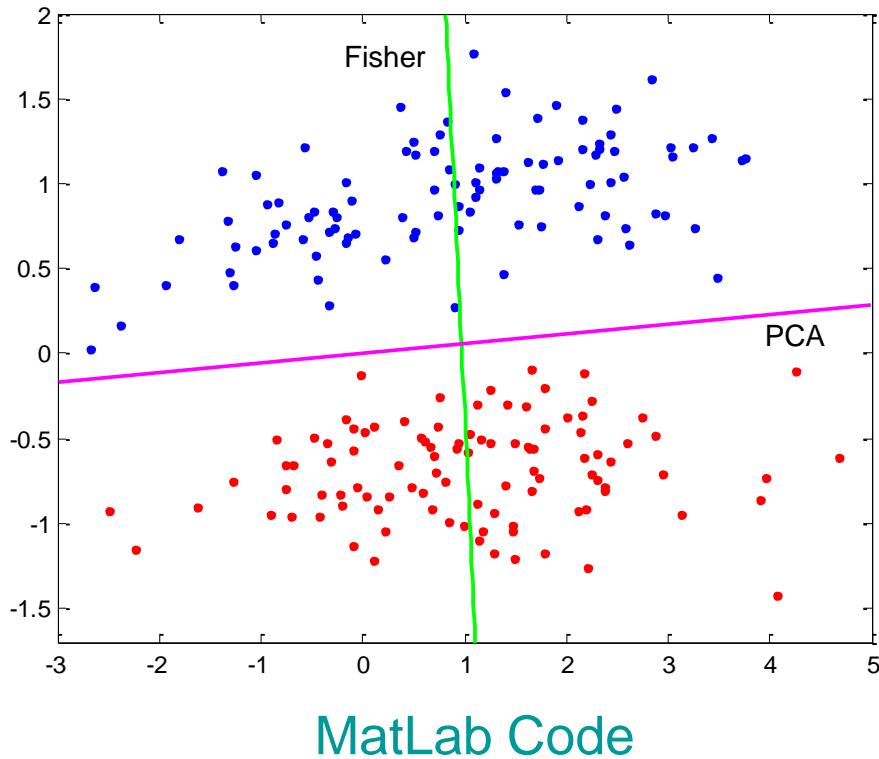
- Visualization of the oil flow data set obtained by projecting the data onto the first two principal components.



[MatLab Code](#)

PCA Vs. Fisher's Discriminant Method

- PCA and the Fisher linear discriminant method can be seen as linear dimensionality reduction techniques.
- PCA is unsupervised and depends only on x_n whereas Fisher linear discriminant also uses class-label information.



- 2D data, 2 classes
- We are interested projecting the data in 1D.
- PCA chooses the direction of maximum variance (magenta curve) leading to strong class overlap.
- The Fisher linear discriminant projects (green curve) to maximize class separation.

PCA For High Dimensional Data

- In many applications of PCA, the number of data points (N) is smaller than the dimensionality (D) of the data space (e.g. 100 images each with 100,000 pixels).
- N points in a D –dimensional space, where $N << D$, defines a linear subspace whose dimensionality is at most $N - 1$.
 - There is no point in applying PCA for values of M that are greater than $N - 1$.
 - If we perform PCA, we will find that at least $D - N + 1$ of the eigenvalues are zero, corresponding to eigenvectors along whose directions the data set has zero variance.

PCA For High Dimensional Data

- Typical algorithms for finding the eigenvectors of a $D \times D$ matrix have a computational cost that scales like $\mathcal{O}(D^3)$.
- Direct application of PCA in e.g. the image example will be computationally infeasible.



PCA For High Dimensional Data

- To resolve the problem, let X to be the $N \times D$ -dimensional centered data matrix, whose n^{th} row is given by $(\mathbf{x}_n - \bar{\mathbf{x}})^T$.

- The covariance $D \times D$ matrix can then be written as:

$$S = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \bar{\mathbf{x}})(\mathbf{x}_n - \bar{\mathbf{x}})^T = \frac{1}{N} X^T X$$

- The corresponding eigenvector equation is:

$$\frac{1}{N} X^T X \mathbf{u}_i = \lambda_i \mathbf{u}_i$$

- Pre-multiplying both sides from the left by X gives:

$$\frac{1}{N} X X^T (\mathbf{X} \mathbf{u}_i) = \lambda_i (\mathbf{X} \mathbf{u}_i)$$



PCA For High Dimensional Data

$$\frac{1}{N} \mathbf{X} \mathbf{X}^T (\mathbf{X} \mathbf{u}_i) = \lambda_i (\mathbf{X} \mathbf{u}_i)$$

- Let us define:

$$\mathbf{v}_i = \mathbf{X} \mathbf{u}_i$$

- The eigenvalue problem becomes:

$$\frac{1}{N} \mathbf{X} \mathbf{X}^T \mathbf{v}_i = \lambda_i \mathbf{v}_i$$

- This is an eigenvector equation for the $N \times N$ matrix $\mathbf{X} \mathbf{X}^T / N$
- This has the same $N - 1$ eigenvalues as the original covariance matrix (which has an additional $D - N + 1$ zero eigenvalues). Note we have $N - 1$ (and not N) eigenvalues because the data are centered ($\mathbf{X} \mathbf{X}^T$ has rank $N - 1$).
- Thus we can solve the eigenvector problem in spaces of lower dimensionality with computational cost $\mathcal{O}(N^3)$.



PCA For High Dimensional Data

- To determine the eigenvectors \mathbf{u}_i , we multiply by \mathbf{X}^T :

$$\frac{1}{N} (\mathbf{X}^T \mathbf{X}) \mathbf{X}^T \mathbf{v}_i = \lambda_i (\mathbf{X}^T \mathbf{v}_i)$$

- The covariance matrix recall has eigenvectors \mathbf{u}_i . Thus with proper rescaling (assuming \mathbf{v}_i is already normalized):

$$\mathbf{u}_i = \frac{1}{(N\lambda_i)^{1/2}} \mathbf{X}^T \mathbf{v}_i$$

- This approach is indeed simple:

- first evaluate $\mathbf{X}\mathbf{X}^T$ and then find its eigenvectors and eigenvalues and
- then compute the eigenvectors in the original data space from the equ. above.



Probabilistic PCA: Maximum Likelihood and EM Algorithm



Probabilistic PCA

PCA as the maximum likelihood solution of a probabilistic latent variable model (Tipping & Bishop 1997, 1999, Roweis, 1998)

- Constrained form of the Gaussian distribution: the number of free parameters is restricted while the model still captures dominant correlations in the data.
- Derive an efficient EM algorithm for PCA.
- Allows dealing with missing values in the dataset.
- Mixture of probabilistic PCA models.

- Tipping, M. E. and C. M. Bishop (1997). Probabilistic principal component analysis. Technical Report NCRG/97/010, Neural Computing Research Group, Aston University.
- Tipping, M. E. and C. M. Bishop (1999b). Probabilistic principal component analysis. *Journal of the Royal Statistical Society, Series B* 21(3), 611–622.
- Roweis, S. (1998). EM algorithms for PCA and SPCA. In M. I. Jordan, M. J. Kearns, and S. A. Solla (Eds.), *Advances in Neural Information Processing Systems*, Volume 10, pp. 626–632. MIT Press.



Probabilistic PCA

PCA as the maximum likelihood solution of a probabilistic latent variable model (Tipping & Bishop [1997](#), [1999](#), Roweis, [1998](#))

- Bayesian treatment of PCA -- the dimensionality of the principal subspace can be found automatically.
- Probabilistic PCA can model class conditional densities and thus be used for classification.
- Can be used generatively: sample from the distribution.

- [Tipping, M. E. and C. M. Bishop \(1999a\). Mixtures of probabilistic principal component analyzers. *Neural Computation* 11\(2\), 443–482.](#)



Probabilistic PCA - Model

- Introduce a latent variable z corresponding to the principal-component subspace.
- Define a Gaussian prior distribution $p(z)$, together with a Gaussian conditional distribution $p(x|z)$ for the observed variable x conditioned on the value of z .
- The prior distribution over z is given by

$$p(z) = \mathcal{N}(z|\mathbf{0}, I)$$

- The conditional distribution of the observed variable x , conditioned on the value of z , is again Gaussian:

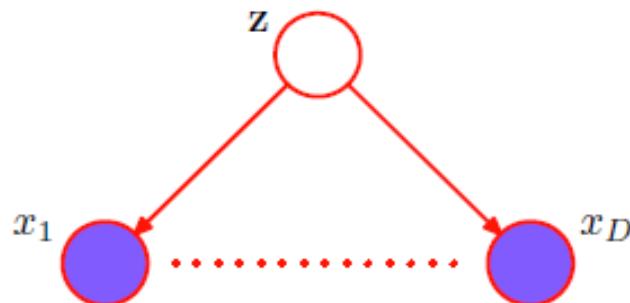
$$p(x|z) = \mathcal{N}(x|Wz + \mu, \sigma^2 I)$$

The mean of x is a linear function of z governed by the $D \times M$ matrix W and the D –dimensional vector μ .



Probabilistic PCA: Naïve Bayes Model

- This model $p(x|z) = \mathcal{N}(x|Wz + \mu, \sigma^2 I)$ is an example of Naive Bayes' model - Conditioned on z , the components of the observed vector $x = (x_1, \dots, x_D)^T$ are assumed to be independent.



- We have 2 parameters in the model:
 - (a) W , its columns span a linear subspace within the data space that corresponds to the principal subspace;
 - (b) σ^2 – variance of the conditional distribution.
- The prior $p(z) = \mathcal{N}(z|\mathbf{0}, I)$ is without loss of generality.

Probabilistic PCA - Generative View Point

- A latent variable model seeks to relate a D -dimensional observation vector x to a corresponding M -dimensional Gaussian latent variable z

$$x = Wz + \mu + \varepsilon$$

where

- z is an M –dimensional Gaussian latent variable
- W is an $(D \times M)$ matrix (the latent space)
- ε is a D –dimensional Gaussian noise
- ε and z are independent
- μ is a parameter vector (non zero mean)

- Factor analysis:

$$\varepsilon \sim \mathcal{N}(\mathbf{0}, \psi)$$

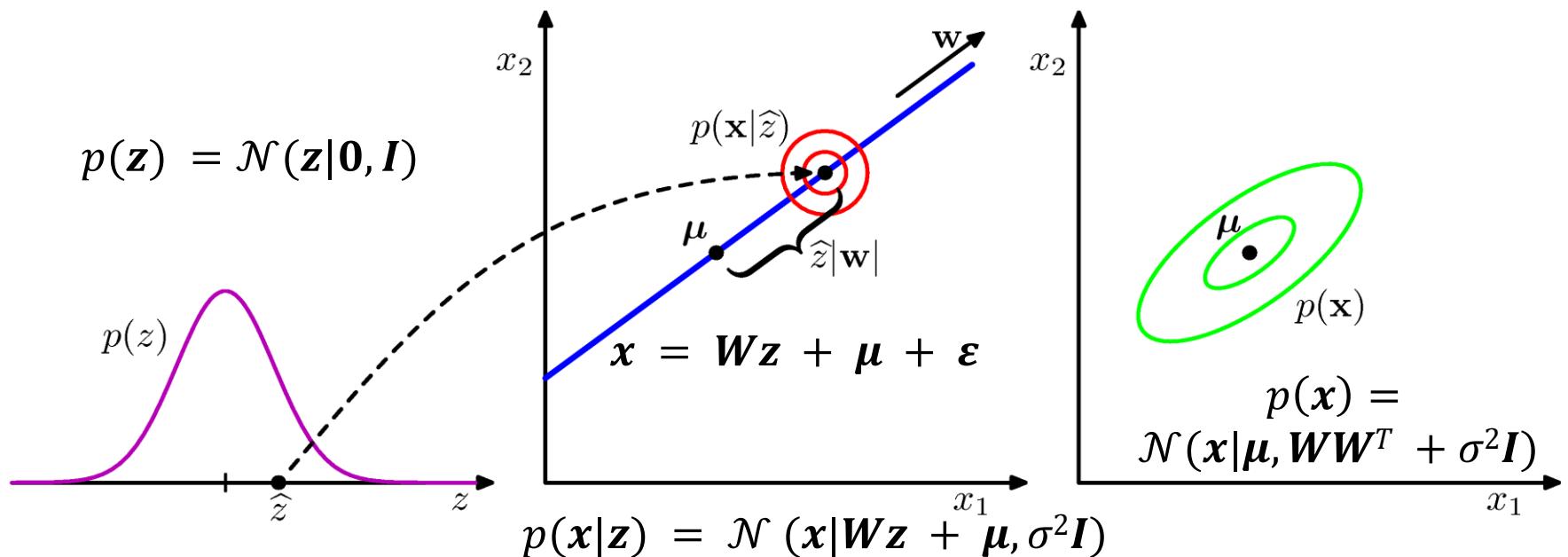
- Probabilistic PCA:

$$\varepsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 I)$$



Probabilistic PCA - Generative View Point

- Mapping from the latent space to the data space.
- Assume here 2D data and 1D latent space.
- An observed x is generated by drawing a value \hat{z} from $p(z)$ & then a value for x from an isotropic Gaussian distribution (red circles) having mean $w\hat{z} + \mu$ and covariance $\sigma^2 I$. The green ellipses are the density contours of $p(x)$.



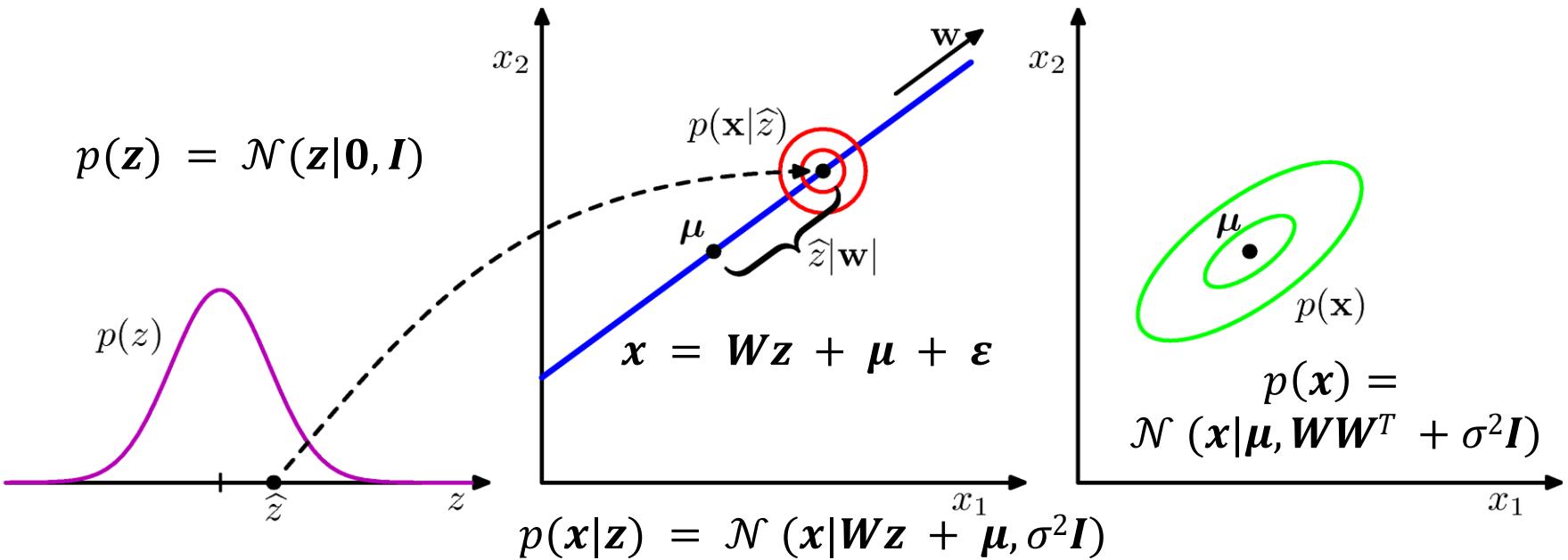
Probabilistic PCA - Generative View Point

- To compute the likelihood function, we need an expression for the marginal distribution $p(x)$ of the observed variable:

$$p(x) = \int p(x | z) p(z) dz$$

- Using Eqs. for linear Gaussian models, this is given as:

$$p(x) = \mathcal{N}(x | \mu, WW^T + \sigma^2 I)$$



Probabilistic PCA - Generative View Point

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \mathbf{W}\mathbf{W}^T + \sigma^2\mathbf{I})$$

- This result can be derived directly by noting that the predictive distribution will be Gaussian and then evaluating its mean and covariance.

$$\mathbb{E}[\mathbf{x}] = \mathbb{E}[\mathbf{Wz} + \boldsymbol{\mu} + \boldsymbol{\varepsilon}] = \mathbf{W}\mathbb{E}[\mathbf{z}] + \boldsymbol{\mu} + \mathbb{E}[\boldsymbol{\varepsilon}] = \boldsymbol{\mu}$$

$$\begin{aligned}\text{var}[\mathbf{x}] &= \mathbb{E}\left[(\mathbf{Wz} + \boldsymbol{\mu} + \boldsymbol{\varepsilon})(\mathbf{Wz} + \boldsymbol{\mu} + \boldsymbol{\varepsilon})^T\right] \\ &= \mathbf{W}\text{var}[\mathbf{z}]\mathbf{W}^T + \mathbb{E}\left[\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^T\right] = \mathbf{W}\mathbf{W}^T + \sigma^2\mathbf{I}\end{aligned}$$

- We used here that \mathbf{z} and $\boldsymbol{\varepsilon}$ are independent, thus uncorrelated.

The Assumed Prior $p(z)$ is Not Restrictive

- The prior we assumed is not restrictive. Indeed consider a more general prior of the form:

$$p(z) = \mathcal{N}(z | \mathbf{m}, \Sigma)$$

- Then the marginal distribution is of the form:

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x} | W\mathbf{m} + \boldsymbol{\mu}, W\Sigma^{-1}W^T + \sigma^2 I)$$

- Thus the marginal distribution has an identical form as before:

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x} | \tilde{\boldsymbol{\mu}}, \widetilde{\mathbf{W}}\widetilde{\mathbf{W}}^T + \sigma^2 I), \quad \begin{aligned}\tilde{\boldsymbol{\mu}} &= W\mathbf{m} + \boldsymbol{\mu} \\ \widetilde{\mathbf{W}} &= \Sigma^{-1/2} W\end{aligned}$$

- We choose the simplest prior: $p(\mathbf{z}) = \mathcal{N}(\mathbf{z} | \mathbf{0}, \mathbf{I})$

Probabilistic PCA - Generative View Point

- There is redundancy in this parametrization.
- There is a whole class of \mathbf{W} 's differing by a rotation of the latent space coordinates that leads to the same predictive distribution.
Indeed:

- Let

$$\widetilde{\mathbf{W}} = \mathbf{W}\mathbf{R}, \text{ where } \mathbf{R} = \text{orthogonal matrix}$$

- Then

$$\widetilde{\mathbf{W}}\widetilde{\mathbf{W}}^T = \mathbf{W}\mathbf{R}\mathbf{R}^T\mathbf{W}^T = \mathbf{W}\mathbf{W}^T$$



Probabilistic PCA - Predictive Distribution

$$p(x) = \mathcal{N}(x | \mu, C) = \mathcal{N}(x | \mu, WW^T + \sigma^2 I)$$

- To compute the predictive distribution, we need to be able to invert C ($D \times D$ matrix). We use the matrix inversion Lemma:

$$C^{-1} = (WW^T + \sigma^2 I)^{-1} = \sigma^{-2}I - \sigma^{-2}WM^{-1}W^T$$

where

$$M = W^T W + \sigma^2 I \text{ } (M \times M \text{ matrix})$$

- The cost of this inversion is reduced from $\mathcal{O}(D^3)$ to $\mathcal{O}(M^3)$!

Probabilistic PCA - Predictive Distribution

$$\mathbf{C}^{-1} = (\mathbf{W}\mathbf{W}^T + \sigma^2\mathbf{I})^{-1} = \sigma^{-2}\mathbf{I} - \sigma^{-2}\mathbf{W}\mathbf{M}^{-1}\mathbf{W}^T$$

$$\mathbf{M} = \mathbf{W}^T\mathbf{W} + \sigma^2\mathbf{I} \text{ } (M \times M \text{ matrix})$$

- To derive this we used the [Woodbury identity](#):

$$(\mathbf{A} + \mathbf{B}\mathbf{D}^{-1}\mathbf{C})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{B}(\mathbf{D} + \mathbf{C}\mathbf{A}^{-1}\mathbf{B})^{-1}\mathbf{C}\mathbf{A}^{-1}$$

- This identity (easy to show with direct substitution) is used often when

- \mathbf{A} is large and diagonal (so easy to invert),
- \mathbf{B} has many rows but few columns $\boxed{\mathbf{B}}$ (and conversely for \mathbf{C} $\boxed{\mathbf{C}}$) so that the rhs is much cheaper to evaluate than the lhs.



Probabilistic PCA - Posterior Distribution

- The posterior distribution can be derived directly from earlier results on linear Gaussian models.
 - We know $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I})$
 - The conditional: $p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}|\mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \sigma^2\mathbf{I})$
 - From these we conclude:

$$p(\mathbf{z} | \mathbf{x}) = \mathcal{N}\left(\mathbf{z} | \mathbf{M}^{-1}\mathbf{W}^T(\mathbf{x} - \boldsymbol{\mu}), \sigma^2\mathbf{M}^{-1}\right), \quad \mathbf{M} = \mathbf{W}^T\mathbf{W} + \sigma^2\mathbf{I}$$

Here we used results for linear Gaussian models:

$$\begin{cases} p(\mathbf{x}) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1}) \\ p(\mathbf{y} | \mathbf{x}) = \mathcal{N}(\mathbf{y} | \mathbf{A}\mathbf{x} + \mathbf{b}, \mathbf{L}^{-1}) \end{cases} \Rightarrow$$

$$p(\mathbf{x} | \mathbf{y}) = \mathcal{N}\left(\mathbf{x} | \boldsymbol{\Sigma}(\boldsymbol{\Lambda}\boldsymbol{\mu} + \mathbf{A}^T\mathbf{L}(\mathbf{y} - \mathbf{b})), \boldsymbol{\Sigma}\right),$$
$$\boldsymbol{\Sigma} = (\boldsymbol{\Lambda} + \mathbf{A}^T\mathbf{L}\mathbf{A})^{-1}$$



Maximum Likelihood PCA

- Consider determining the model parameters using maximum likelihood. Using

$$p(\mathbf{x}) = \mathcal{N}\left(\mathbf{x} \mid \boldsymbol{\mu}, \underbrace{\mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I}}_C\right)$$

we derive:

$$\begin{aligned} \ln p(\mathbf{X} \mid \boldsymbol{\mu}, \mathbf{W}, \sigma^2) &= \sum_{n=1}^N \ln p(\mathbf{x}_n \mid \boldsymbol{\mu}, \mathbf{W}, \sigma^2) \\ &= -\frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln |\mathbf{C}| - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^T \mathbf{C}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) \end{aligned}$$

- Setting the derivative wrt $\boldsymbol{\mu}$ equal to zero gives:

$$\boldsymbol{\mu} = \sum_{n=1}^N \mathbf{x}_n \Big/ N \equiv \bar{\mathbf{x}}$$

Maximum Likelihood PCA

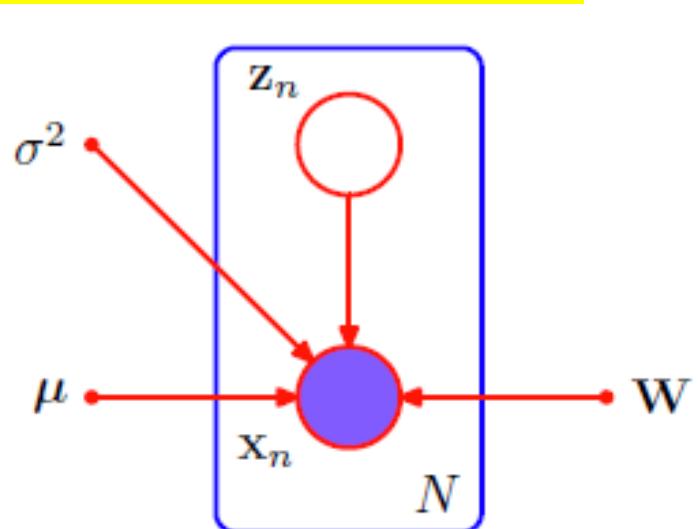
- The log-likelihood is then simplified as:

$$\ln p(X | \mu, W, \sigma^2) = -\frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln|C| - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \bar{\mathbf{x}})^T C^{-1} (\mathbf{x}_n - \bar{\mathbf{x}})$$

or

$$\ln p(X | \mu, W, \sigma^2) = -\frac{N}{2} \left\{ D \ln(2\pi) + \ln|C| + \text{Tr}(C^{-1} S) \right\},$$

$$S = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \bar{\mathbf{x}})(\mathbf{x}_n - \bar{\mathbf{x}})^T$$



- Maximization wrt W and σ^2 can also be done analytically:

$$W_{ML} = U_M \left(L_M - \sigma^2 I \right)^{1/2} R$$

- Tipping, M. E. and C. M. Bishop (1999b). Probabilistic principal component analysis. *Journal of the Royal Statistical Society, Series B* 21(3), 611–622.
- Roweis, S. (1998). EM algorithms for PCA and SPCA. In M. I. Jordan, M. J. Kearns, and S. A. Solla (Eds.), *Advances in Neural Information Processing Systems, Volume 10*, pp. 626–632. MIT Press.

Maximum Likelihood PCA

$$\mathbf{W}_{ML} = \mathbf{U}_M \left(\mathbf{L}_M - \sigma^2 \mathbf{I} \right)^{1/2} \mathbf{R}$$

- \mathbf{U}_M is a $D \times M$ matrix whose columns are given by any subset (of size M) of the eigenvectors of the data covariance matrix \mathbf{S} .
- \mathbf{L}_M is the $M \times M$ diagonal matrix with elements given by the corresponding eigenvalues λ_i of \mathbf{S} .
- \mathbf{R} is an arbitrary $M \times M$ orthogonal matrix.
- The max of the likelihood function is obtained when the M eigenvectors are chosen to be those whose eigenvalues are the M largest (all other solutions being saddle points).
- For eigenvectors arranged in order of decreasing λ_i values, the M principal eigenvectors are $\mathbf{u}_1, \dots, \mathbf{u}_M$. The columns of \mathbf{W} then define the principal subspace as in standard PCA.

- Tipping, M. E. and C. M. Bishop (1999b). Probabilistic principal component analysis. *Journal of the Royal Statistical Society, Series B* **21**(3), 611–622.



Maximum Likelihood PCA

- The corresponding MLE solution for σ^2 is given as:

$$\sigma_{ML}^2 = \frac{1}{D-M} \sum_{i=M+1}^N \lambda_i$$

i.e. the average of the discarded eigenvalues.

- R in $W_{ML} = U_M (L_M - \sigma^2 I)^{1/2} R$ is a rotation matrix in the M dimensional latent space.
- Substituting this into the predictive variance $C = W_{ML} W_{ML}^T + \sigma_{ML}^2 I$ and using $R^T R = I$, we see that C is independent of R .
- The predictive density is unchanged by rotations in the latent space (statistical non-identifiability).



Maximum Likelihood PCA

$$\mathbf{W}_{ML} = \mathbf{U}_M \left(\mathbf{L}_M - \sigma^2 \mathbf{I} \right)^{1/2} \mathbf{R}$$

- For $\mathbf{R} = \mathbf{I}$, the columns of \mathbf{W} are the principal component eigenvectors scaled by $(\lambda_i - \sigma^2)^{1/2}$.
- Interpretation: for the convolution of independent Gaussian distributions, here of \mathbf{z} and the noise $\boldsymbol{\varepsilon}$, the variances are additive.

$$\mathbf{x} = \mathbf{W}\mathbf{z} + \boldsymbol{\mu} + \boldsymbol{\varepsilon}$$

- The variance λ_i in the direction of an eigenvector \mathbf{u}_i is composed of the sum of a contribution $\lambda_i - \sigma^2$ from the projection of the unit-variance latent space $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I})$ into data space through the corresponding column of \mathbf{W} , plus an isotropic contribution of variance σ^2 added in all directions by the noise model $\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$.



Maximum Likelihood PCA

- Consider the variance of the predictive distribution $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\mathbf{m}, \mathbf{C})$, $\mathbf{C} = \mathbf{W}\mathbf{W}^T + \sigma^2\mathbf{I}$ along some direction defined by the unit vector $\boldsymbol{\nu}$ given by $\boldsymbol{\nu}^T \mathbf{C} \boldsymbol{\nu}$.
- Let $\boldsymbol{\nu}$ be orthogonal to the principal subspace \mathbf{U} , i.e. $\boldsymbol{\nu}$ is given by some linear combination of the discarded eigenvectors.
- Then $\boldsymbol{\nu}^T \mathbf{U} = 0$ and hence

$$\boldsymbol{\nu}^T \mathbf{C} \boldsymbol{\nu} = \boldsymbol{\nu}^T (\mathbf{W}\mathbf{W}^T + \sigma^2\mathbf{I}) \boldsymbol{\nu} = \sigma^2$$

- Thus the model predicts a noise variance orthogonal to the principal subspace which it was shown to be the average of the discarded eigenvalues.

$$\sigma_{ML}^2 = \frac{1}{D-M} \sum_{i=M+1}^N \lambda_i$$



Maximum Likelihood PCA

- Consider the variance of the predictive distribution $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\mathbf{m}, \mathbf{C})$, $\mathbf{C} = \mathbf{W}\mathbf{W}^T + \sigma^2\mathbf{I}$ along some direction defined by the unit vector \mathbf{v} given by $\mathbf{v}^T \mathbf{C} \mathbf{v}$.
- Let now $\mathbf{v} = \mathbf{u}_i$ where \mathbf{u}_i is one of the retained eigenvectors defining the principal subspace.
- Then using $\mathbf{W}_{ML} = \mathbf{U}_M (\mathbf{L}_M - \sigma^2 \mathbf{I})^{1/2} \mathbf{R}$, and $\mathbf{u}_i^T \mathbf{u}_j = 0, j = 1, \dots, M$ we see that:

$$\mathbf{v}^T \mathbf{C} \mathbf{v} = \mathbf{u}_i^T (\mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I}) \mathbf{u}_i = (\lambda_i - \sigma^2) + \sigma^2 = \lambda_i$$

i.e. the model correctly captures the variance of the data along the principal axes.

- As shown earlier, the variance in all remaining directions is approximated with a single value of σ^2 . Variance is ‘lost’ in the projections.



Maximum Likelihood PCA

$$\mathbf{W}_{ML} = \mathbf{U}_M \left(\mathbf{L}_M - \sigma^2 \mathbf{I} \right)^{1/2} \mathbf{R}$$

$$\sigma_{ML}^2 = \frac{1}{D-M} \sum_{i=M+1}^N \lambda_i$$

- We construct the MLE model by solving the underlying eigenvalue problem for the data covariance matrix and evaluate \mathbf{W} and σ^2 from the Eqs. above. We choose $\mathbf{R} = \mathbf{I}$.
- Note: if the MLE model is found using optimization methods or via the EM algorithm (see following slides), \mathbf{R} is arbitrary and the columns of \mathbf{W} not orthogonal.

In this case, orthogonality can be enforced

- as postprocessing or
 - by modifying the EM algorithm.
-
- Ahn, J. H. and J. H. Oh (2003). [A constrained EM algorithm for principal component analysis](#). *Neural Computation* 15(1), 57–65.



Maximum Likelihood PCA

- If we consider no dimensionality reduction ($M = D$), then

$$\mathbf{U}_M = \mathbf{U}, \quad \mathbf{L}_M = \mathbf{L},$$

and using $\mathbf{U}\mathbf{U}^T = \mathbf{I}$, $\mathbf{R}\mathbf{R}^T = \mathbf{I}$, we see that the covariance \mathbf{C} of the marginal distribution for \mathbf{x} becomes

$$\mathbf{C} = \mathbf{W}\mathbf{W}^T + \sigma^2\mathbf{I} = \mathbf{U}\left(\mathbf{L} - \sigma^2\mathbf{I}\right)^{1/2}\mathbf{R}\mathbf{R}^T\left(\mathbf{L} - \sigma^2\mathbf{I}\right)^{1/2}\mathbf{U}^T + \sigma^2\mathbf{I} = \mathbf{U}\mathbf{L}\mathbf{U}^T = \mathbf{S}$$

- This result is the standard MLE solution for an unconstrained Gaussian distribution in which the covariance matrix is given by the sample covariance.

Efficient Evaluation of PPCA Density

- Since \mathbf{C} is not full rank, we use the matrix inversion lemma:

$$\mathbf{C}^{-1} = \frac{1}{\sigma^2} \left[\mathbf{I} - \mathbf{W} \left(\mathbf{W}^T \mathbf{W} + \sigma^2 \mathbf{I} \right)^{-1} \mathbf{W}^T \right]$$

- Plugging the MLE estimates for \mathbf{W}, σ^2 gives: $\mathbf{W} = \mathbf{U}_M \left(\mathbf{L}_M - \sigma^2 \mathbf{I} \right)^{1/2}$

$$\mathbf{C}^{-1} = \frac{1}{\sigma^2} \left[\mathbf{I} - \mathbf{U}_M \left(\mathbf{L}_M - \sigma^2 \mathbf{I} \right)^{1/2} \left(\mathbf{L}_M - \sigma^2 \mathbf{I} + \sigma^2 \mathbf{I} \right)^{-1} \left(\mathbf{L}_M - \sigma^2 \mathbf{I} \right)^{1/2} \mathbf{U}_M^T \right]$$

$$= \frac{1}{\sigma^2} \left[\mathbf{I} - \mathbf{U}_M \left(\mathbf{L}_M - \sigma^2 \mathbf{I} \right)^{1/2} \mathbf{L}_M^{-1} \left(\mathbf{L}_M - \sigma^2 \mathbf{I} \right)^{1/2} \mathbf{U}_M^T \right]$$

$$= \frac{1}{\sigma^2} \left[\mathbf{I} - \mathbf{U}_M \text{diag} \left(1 - \frac{\sigma^2}{\lambda_j} \right) \mathbf{U}_M^T \right]$$

- Similarly: $\log |\mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I}| = (D - M) \log \sigma^2 + \sum_{i=1}^M \log \lambda_i$

Use the matrix inversion lemma $\det(\mathbf{A} + \mathbf{U}\mathbf{W}\mathbf{V}^T) = \det(\mathbf{W}^{-1} + \mathbf{V}^T \mathbf{A}^{-1} \mathbf{U}) \det \mathbf{W} \det \mathbf{A}$

Efficient Evaluation of PPCA Density

- Since C is not full rank, we can use the matrix inversion lemma:

$$|C| = |WW^T + \sigma^2 I_{D \times D}| = |U_M (L_M - \sigma^2 I_{M \times M}) U^T + \sigma^2 I_{D \times D}| \Rightarrow$$

$$|C| = \left\| (L_M - \sigma^2 I_{M \times M})^{-1} + U^T \sigma^{-2} I_{M \times M} U \right\| \left\| (L_M - \sigma^2 I_{M \times M}) \right\| |\sigma^2 I_{D \times D}|$$

- Here we used the MLE value for $W = U_M (L_M - \sigma^2 I)^{1/2}$ and the matrix inversion lemma $\det(A + UWV^T) = \det(W^{-1} + V^T A^{-1} U) \det W \det A$
- Taking logs and evaluating the above determinants gives:

$$\log |C| = \left(\sum_{i=1}^M \log \lambda_i - M \log \sigma^2 - \sum_{i=1}^M \log(\lambda_i - \sigma^2) \right) + \sum_{i=1}^M \log(\lambda_i - \sigma^2) + D \log \sigma^2$$

$$\log |WW^T + \sigma^2 I| = (D - M) \log \sigma^2 + \sum_{i=1}^M \log \lambda_i$$

Maximum Likelihood PCA

- PCA is generally expressed as a projection of points from the D -dimensional dataspace onto an M -dimensional subspace.
- Probabilistic PCA is seen as the mapping $\mathbf{z} \rightarrow \mathbf{x}, \mathbf{x} = \mathbf{W}\mathbf{z} + \boldsymbol{\mu} + \boldsymbol{\varepsilon}$.
- The reverse mapping $\mathbf{x} \rightarrow \mathbf{z}$ is computed using the posterior:

$$p(\mathbf{z} | \mathbf{x}) = \mathcal{N}\left(\mathbf{z} | \mathbf{M}^{-1}\mathbf{W}^T(\mathbf{x} - \boldsymbol{\mu}), \sigma^2\mathbf{M}^{-1}\right), \quad \mathbf{M} = \mathbf{W}^T\mathbf{W} + \sigma^2\mathbf{I}$$

- Every point in data space is characterized by its posterior mean and covariance in latent space:

$$\mathbb{E}[z | x] = \mathbf{M}^{-1}\mathbf{W}_{ML}^T(\mathbf{x} - \bar{\mathbf{x}}), \quad \text{var}[z | x] = \sigma^2\mathbf{M}^{-1}$$

- This projects to a point in data space given by:

$$\mathbf{W}\mathbb{E}[z | x] + \boldsymbol{\mu}$$

- This is identical form as regularized linear regression.



Maximum Likelihood PCA vs Standard PCA

- In the limit $\sigma^2 \rightarrow 0$, the posterior mean becomes:

$$\mathbb{E}[z | x] = M^{-1} W_{ML}^T (x - \bar{x}), M = W_{ML}^T W_{ML} + \sigma^2 I \rightarrow W_{ML}^T W_{ML} \Rightarrow$$

$$\mathbb{E}[z | x] = (W_{ML}^T W_{ML})^{-1} W_{ML}^T (x - \bar{x})$$

- Now substitute $W_{ML} = U_M (L_M - \sigma^2 I)^{1/2} R = U_M L_M^{1/2}$ for $\sigma^2 \rightarrow 0$ with $R = I$ (for consistency with PCA). Then

$$\mathbb{E}[z | x] = (W_{ML}^T W_{ML})^{-1} W_{ML}^T (x - \bar{x}) = L_M^{-1/2} U_M^T (x - \bar{x})$$

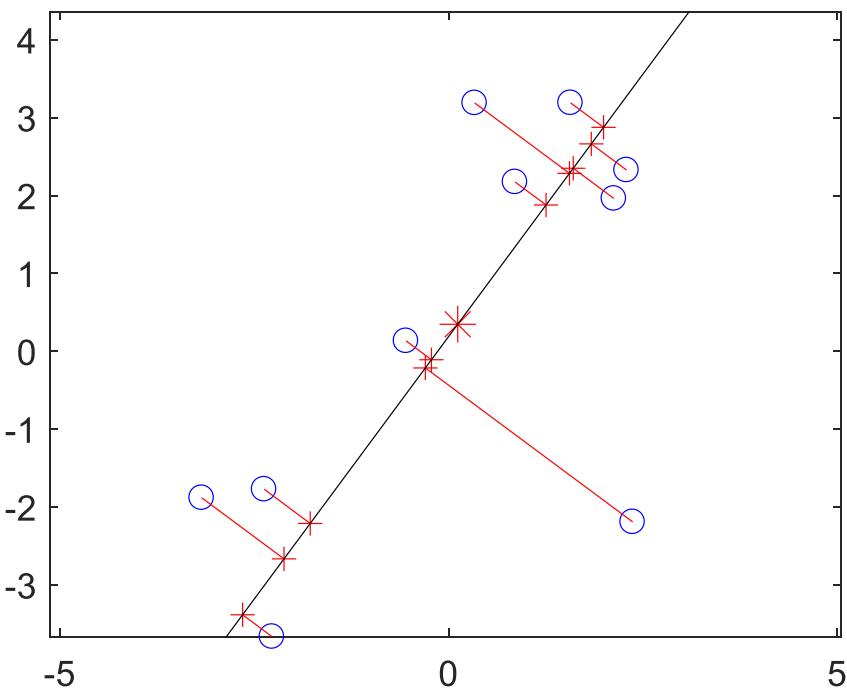
- This is an orthogonal projection of the data point onto the latent space, i.e. for the limit $\sigma^2 \rightarrow 0$, we recover the standard PCA model. The posterior covariance

$$\text{var}[z | x] = \sigma^2 M^{-1} \rightarrow 0 \text{ and the density becomes singular.}$$

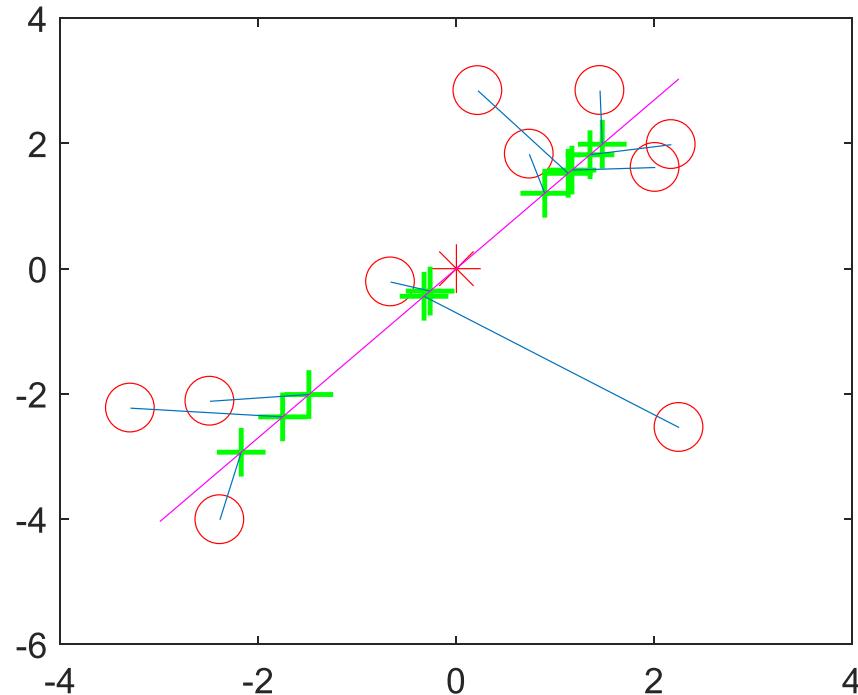
- For $\sigma^2 > 0$, the latent projection is shifted towards the origin, relative to the orthogonal projection.

PPCA Versus PCA

- Consider PCA and PPCA where $D = 2$ and $M = 1$. The red star is the data mean. In PCA the points are orthogonally projected onto the line. In PPCA the projection is no longer orthogonal and the reconstructions are shrunk towards the data mean (red star).



Run [pcaDemo2d](#)
From [PMTK3](#)



Run [ppcaDemo2d](#)
From [PMTK3](#)



Probabilistic PCA: Number of DOF

- PPCA defines a multivariate Gaussian distribution in which the number of DOF (independent parameters) can be controlled while still capturing the dominant correlations in the data.
- A general Gaussian distribution has $D(D + 1)/2$ independent parameters in its covariance matrix and D parameters in its mean.
 - The number of parameters scales as D^2 .
- For a diagonal covariance matrix we have D independent parameters and the number of parameters grows linearly. However, the variables are independent and hence this model cannot express any correlations between them.

Probabilistic PCA: Number of DOF

- In PPCA, the M most significant correlations are captured while the total number of parameters grows linearly with D .
- We can see this by evaluating the DOF in PPCA:
 - The covariance \mathbf{C} depends on \mathbf{W} ($D \times M$), and σ^2 : $\mathbf{C} = \mathbf{W}\mathbf{W}^T + \sigma^2\mathbf{I}$ total parameters $DM + 1$.
 - We need to subtract the redundancy associated with rotations of the coordinate system in the latent space.
 - \mathbf{R} is $M \times M$. In the 1st column there are $M - 1$ independent parameters (must be normalized). In the 2nd column there are $M - 2$ independent parameters (normalized & orthogonal to the 1st column), etc. \mathbf{R} has a total of $M(M - 1)/2$ independent parameters.
- The number of degrees of freedom in \mathbf{C} grows linearly with D

$$D \times M + 1 - M \times (M - 1)/2$$



Probabilistic PCA: Number of DOF

- The number of DOF in \mathcal{C} grows linearly in D for a given M :

$$D \times M + 1 - M \times (M - 1)/2$$

- $M = D - 1$: recover the standard result for a full covariance Gaussian.
 - The variance along $D - 1$ linearly independent directions is controlled by the columns of \mathbf{W} , and
 - The variance along the remaining direction is given by σ^2 .
- $M = 0$: equivalent to the isotropic covariance case.

EM Algorithm for PCA



EM Algorithm for PCA

- The PPCA model involves a marginalization over a continuous latent space \mathbf{z} . For each x_n there is a corresponding \mathbf{z}_n .
- Use EM to find the MLE of the model parameters.
- For high D , there are advantages using iteratively EM vs. working directly with the sample covariance.
- This approach can be extended to factor analysis for which there is no closed-form solution.
- Can also be used when values are missing & for mixture models.
- EM requires the complete-data log likelihood function:

$$\ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\mu}, \mathbf{W}, \sigma^2) = \sum_{n=1}^N \left\{ \ln p(x_n | z_n) + \ln p(z_n) \right\}$$

- In the following, we substitute $\boldsymbol{\mu}$ with the sample mean $\bar{\mathbf{x}}$.

EM Algorithm for PCA

- We first take the expectation of the complete-data log likelihood with respect to the posterior distribution of the latent distribution evaluated using ‘old’ parameter values.
- Maximization of this expected complete data log likelihood then yields the ‘new’ parameter values.

$$\ln p(X, Z | \mu, W, \sigma^2) = \sum_{n=1}^N \{ \ln p(x_n | z_n) + \ln p(z_n) \}$$

- The n^{th} row of Z is given by z_n .
- Recall that $p(z) = \mathcal{N}(z|\mathbf{0}, I)$, $p(x|z) = \mathcal{N}(x|Wz + \mu, \sigma^2 I)$
- We can now write the expectation with respect to the posterior distribution over the latent variables.



EM Algorithm for PCA

$$\mathbb{E} \left[\ln p(X, Z | \mu, W, \sigma^2) \right] = -\sum_{n=1}^N \left\{ \begin{aligned} & \frac{D}{2} \ln(2\pi\sigma^2) + \frac{1}{2} \text{Tr}(\mathbb{E}[z_n z_n^T]) \\ & + \frac{1}{2\sigma^2} \|x_n - \mu\|^2 - \frac{1}{\sigma^2} \mathbb{E}[z_n]^T W^T (x_n - \mu) \\ & + \frac{1}{2\sigma^2} \text{Tr}(\mathbb{E}[z_n z_n^T] W^T W) + \frac{M}{2 \ln(2\pi)} \end{aligned} \right\}$$

□ **E-Step:** We use the old parameters to evaluate:

$$\mathbb{E}[z_n] = M^{-1}W^T(x_n - \bar{x}) \quad \mathbb{E}[z_n z_n^T] = \sigma^2 M^{-1} + \mathbb{E}[z_n] \mathbb{E}[z_n]^T$$

□ This follows directly from

$$p(z | x) = \mathcal{N}(z | M^{-1}W^T(x - \mu), \sigma^2 M^{-1}), \quad M = W^T W + \sigma^2 I$$

together with the standard result

$$\mathbb{E}[z_n z_n^T] = \text{cov}[z_n] + \mathbb{E}[z_n] \mathbb{E}[z_n]^T$$

EM Algorithm for PCA

$$\mathbb{E} \left[\ln p(X, Z | \mu, W, \sigma^2) \right] = -\sum_{n=1}^N \left\{ \begin{aligned} & \frac{D}{2} \ln(2\pi\sigma^2) + \frac{1}{2} \text{Tr}(\mathbb{E}[z_n z_n^T]) \\ & + \frac{1}{2\sigma^2} \|x_n - \mu\|^2 - \frac{1}{\sigma^2} \mathbb{E}[z_n]^T W^T (x_n - \mu) \\ & + \frac{1}{2\sigma^2} \text{Tr}(\mathbb{E}[z_n z_n^T] W^T W) + \frac{M}{2 \ln(2\pi)} \end{aligned} \right\}$$

□ **M-Step:** We maximize with respect to W and σ^2 keeping the posterior statistics fixed. We obtain:

$$\begin{aligned} W_{new} &= \left[\sum_{n=1}^N (x_n - \bar{x}) \mathbb{E}[z_n]^T \right] \left[\sum_{n=1}^N \mathbb{E}[z_n z_n^T] \right]^{-1} \\ \sigma_{new}^2 &= \frac{1}{ND} \sum_{n=1}^N \left\{ \begin{aligned} & \|x_n - \bar{x}\|^2 - 2 \mathbb{E}[z_n]^T W_{new} (x_n - \bar{x}) \\ & + \text{Tr}(\mathbb{E}[z_n z_n^T] W_{new}^T W_{new}) \end{aligned} \right\} \end{aligned}$$

Proof of the M-Step Equations

$$\mathbb{E} \left[\ln p(X, Z | \mu, W, \sigma^2) \right] = -\sum_{n=1}^N \left\{ \begin{aligned} & \frac{D}{2} \ln(2\pi\sigma^2) + \frac{1}{2} \text{Tr} \left(\mathbb{E} [z_n z_n^T] \right) \\ & + \frac{1}{2\sigma^2} \|x_n - \mu\|^2 - \frac{1}{\sigma^2} \mathbb{E} [z_n]^T W^T (x_n - \mu) \\ & + \frac{1}{2\sigma^2} \text{Tr} \left(\mathbb{E} [z_n z_n^T] W^T W \right) + \frac{M}{2 \ln(2\pi)} \end{aligned} \right\}$$

- The two M-equations are derived by setting the derivatives wrt W and σ^2 equal to zero:

$$\frac{\partial \mathbb{E} \left[\ln p(X, Z | \mu, W, \sigma^2) \right]}{\partial W} = \sum_{n=1}^N \left\{ \frac{1}{\sigma^2} (x_n - \mu) \mathbb{E} [z_n]^T - \frac{1}{\sigma^2} W \mathbb{E} [z_n z_n^T] \right\} = 0$$

$$\frac{\partial \mathbb{E} \left[\ln p(X, Z | \mu, W, \sigma^2) \right]}{\partial \sigma^2} = \sum_{n=1}^N \left\{ \begin{aligned} & -\frac{D}{2\sigma^2} - \frac{1}{\sigma^4} \mathbb{E} [z_n]^T W^T (x_n - \mu) + \frac{1}{2\sigma^4} \|x_n - \mu\|^2 \\ & + \frac{1}{2\sigma^4} \text{Tr} \left(\mathbb{E} [z_n z_n^T] W^T W \right) \end{aligned} \right\} = 0$$

- Here, we used $\frac{\partial}{\partial A} \text{Tr}(ABA^T) = A(B + B^T)$, $\frac{\partial}{\partial A} \text{Tr}(AB) = B^T$

EM Algorithm for PCA

- Initialize the parameters
- Compute the sufficient statistics of the latent space posterior distribution in the E-Step

$$\mathbb{E}[\mathbf{z}_n] = \mathbf{M}^{-1} \mathbf{W}^T (\mathbf{x}_n - \bar{\mathbf{x}}) \quad \mathbb{E}[\mathbf{z}_n \mathbf{z}_n^T] = \sigma^2 \mathbf{M}^{-1} + \mathbb{E}[\mathbf{z}_n] \mathbb{E}[\mathbf{z}_n^T]$$

- Revise the parameter values in the M-Step.

$$\mathbf{W}_{new} = \left[\sum_{n=1}^N (\mathbf{x}_n - \bar{\mathbf{x}}) \mathbb{E}[\mathbf{z}_n]^T \right] \left[\sum_{n=1}^N \mathbb{E}[\mathbf{z}_n \mathbf{z}_n^T] \right]^{-1}$$

$$\sigma_{new}^2 = \frac{1}{ND} \sum_{n=1}^N \left\{ \|\mathbf{x}_n - \bar{\mathbf{x}}\|^2 - 2 \left[\mathbf{z}_n \right]^T \mathbf{W}_{new}^T (\mathbf{x}_n - \bar{\mathbf{x}}) \right\} \\ + Tr \left(\mathbb{E}[\mathbf{z}_n \mathbf{z}_n^T] \mathbf{W}_{new}^T \mathbf{W}_{new} \right)$$

EM Algorithm for PCA: Computational Cost

- Each cycle of the EM algorithm can be computationally more efficient than conventional PCA in high dimensions.
- The eigen-decomposition of the covariance matrix requires $\mathcal{O}(D^3)$ computation. If interested only in the first M eigenvectors, we can use algorithms that are $\mathcal{O}(MD^2)$.
- However, the evaluation of the covariance matrix itself

$$\mathbf{S} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \bar{\mathbf{x}})(\mathbf{x}_n - \bar{\mathbf{x}})^T$$

takes $\mathcal{O}(ND^2)$ computations.

- Algorithms such as the snapshot method ([Sirovich, 1987](#)), assume that the eigenvectors are linear combinations of the data vectors and avoid direct evaluation of the covariance matrix but are $\mathcal{O}(N^3)$ and hence unsuited to large data sets.

- Sirovich, L. (1987). [Turbulence and the dynamics of coherent structures](#). *Quarterly Applied Mathematics* **45**(3), 561–590.



EM Algorithm for PCA: Computational Cost

- The EM algorithm described here does not construct the covariance matrix explicitly.
- Instead, the most computationally demanding steps are those involving sums over the data set that are $\mathcal{O}(NDM)$.

$$\mathbf{W}_{new} = \left[\sum_{n=1}^N (\mathbf{x}_n - \bar{\mathbf{x}}) \mathbb{E}[\mathbf{z}_n]^T \right] \left[\sum_{n=1}^N \mathbb{E}[\mathbf{z}_n \mathbf{z}_n^T] \right]^{-1}$$

- For large D , and $M \ll D$, $\mathcal{O}(NDM)$ is much better compared to $\mathcal{O}(ND^2)$ (that is needed to compute the sample covariance)

This offsets the iterative nature of the EM algorithm.



EM Algorithm for PCA: Online form

- The EM algorithm can be implemented in an on-line form: each D -dim data point is read/processed/then discarded before the next data point is considered.
- To see this, note that the quantities evaluated in the E-Step (an M -dimensional vector and an $M \times M$ matrix) can be computed for each data point separately.

$$\mathbb{E}[z_n] = M^{-1}W^T(x_n - \bar{x}) \quad \mathbb{E}[z_n z_n^T] = \sigma^2 M^{-1} + \mathbb{E}[z_n] \mathbb{E}[z_n]^T$$

- In the M – Step, we accumulate the sums over data points incrementally - advantageous if both N & D are large.

$$W_{new} = \left[\sum_{n=1}^N (x_n - \bar{x}) \mathbb{E}[z_n]^T \right] \left[\sum_{n=1}^N \mathbb{E}[z_n z_n^T] \right]^{-1}$$
$$\sigma_{new}^2 = \frac{1}{ND} \sum_{n=1}^N \left\{ \begin{aligned} & \left\| x_n - \bar{x} \right\|^2 - 2 \mathbb{E}[z_n]^T W_{new}^T (x_n - \bar{x}) \\ & + Tr(\mathbb{E}[z_n z_n^T] W_{new}^T W_{new}) \end{aligned} \right\}$$



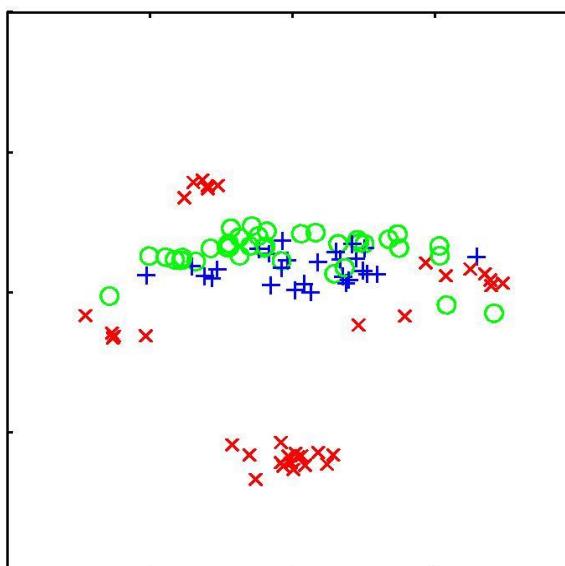
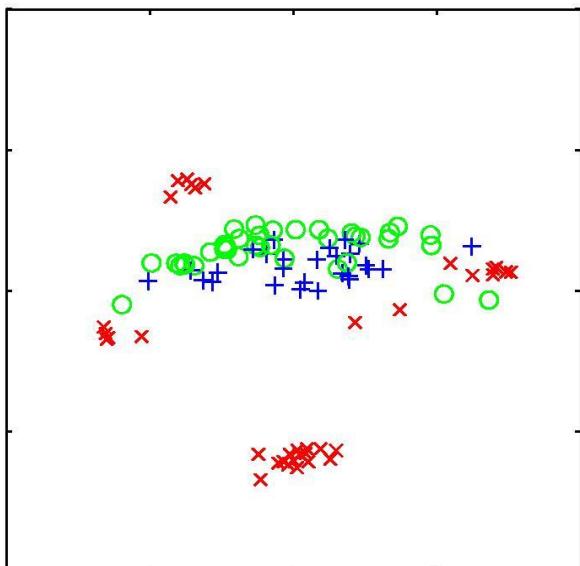
EM Algorithm for PCA: Missing Values

- ❑ Because we now have a fully probabilistic model for PCA, we can **deal with missing data**, provided that it is **missing at random**, by marginalizing over the distribution of the unobserved variables.
- ❑ Again these missing values can be treated using the EM algorithm.
- ❑ We give an example of the use of this approach for data visualization in an example next.



EM Algorithm for PCA: Missing Values

- PPCA visualization for the first 100 data points of the oil flow data set.
- Left: the posterior mean projections of the data points on the principal subspace.
- Right: randomly omitting 30% of the variable values and using EM to handle the missing values. Even though each data point has at least one missing measurement, the plot is similar to the one obtained without missing values.



[Matlab Implementation](#)

EM Algorithm for PCA: Limit $\sigma^2 \rightarrow 0$

- When $\sigma^2 \rightarrow 0$, EM corresponds to standard PCA ([Roweis, 1998](#))
- Defining $\tilde{\mathbf{X}}$ a matrix of size $N \times D$ whose n^{th} row is given by $\mathbf{x}_n - \bar{\mathbf{x}}$.
- Defining a matrix Ω of size $M \times N$ whose n^{th} column is given by the vector $\mathbb{E}[\mathbf{z}_n]$.
- The E-Step for $\sigma^2 \rightarrow 0$ becomes

$$\mathbb{E}[\mathbf{z}_n] = \mathbf{M}^{-1} \mathbf{W}^T (\mathbf{x}_n - \bar{\mathbf{x}}) \Rightarrow \Omega = (\mathbf{W}_{old}^T \mathbf{W}_{old})^{-1} \mathbf{W}_{old}^T \tilde{\mathbf{X}}^T$$

This is simply the orthogonal projection of the data points on the current estimate for the principal subspace.

- Roweis, S. (1998). [EM algorithms for PCA and SPCA](#). In M. I. Jordan, M. J. Kearns, and S. A. Solla (Eds.), [Advances in Neural Information Processing Systems, Volume 10](#), pp. 626–632. MIT Press.



EM Algorithm for PCA: Limit $\sigma^2 \rightarrow 0$

- When $\sigma^2 \rightarrow 0$, EM corresponds to standard PCA
- $\tilde{\mathbf{X}}$ an $N \times D$ matrix whose n^{th} row is given by $\mathbf{x}_n - \bar{\mathbf{x}}$.
- Ω a matrix $M \times N$ whose n^{th} column is given by $\mathbb{E}[\mathbf{z}_n]$.
- Noting that

$$\sigma^2 \rightarrow 0 \Rightarrow \mathbb{E}[\mathbf{z}_n \mathbf{z}_n^T] = \sigma^2 \mathbf{M}^{-1} + \mathbb{E}[\mathbf{z}_n] \mathbb{E}[\mathbf{z}_n]^T \rightarrow \mathbb{E}[\mathbf{z}_n] \mathbb{E}[\mathbf{z}_n]^T$$

the M –Step takes the form:

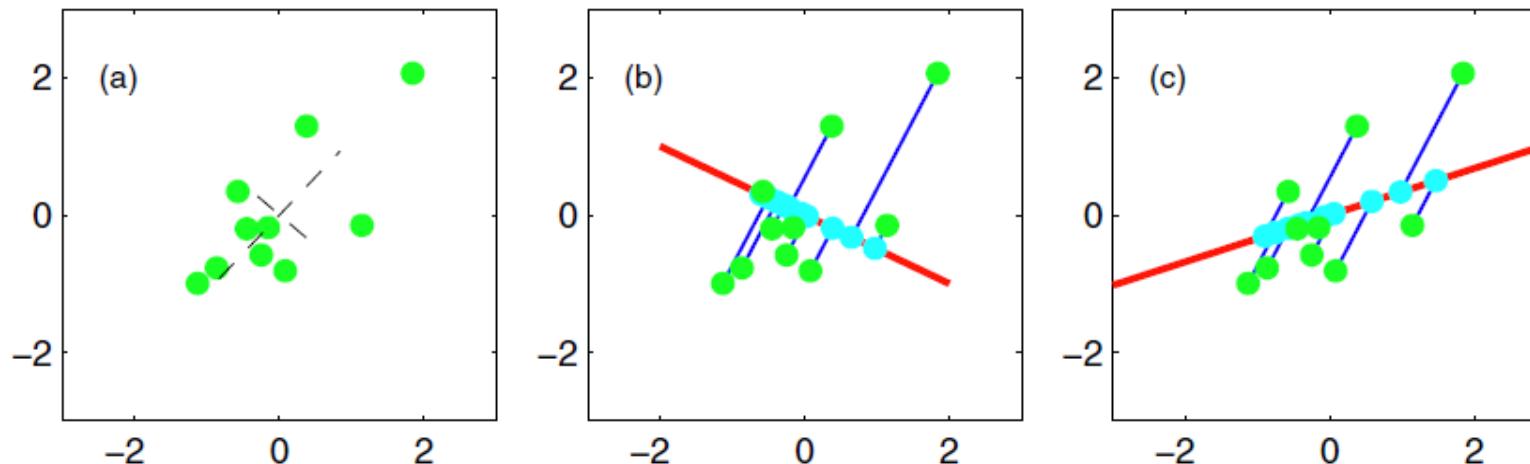
$$\mathbf{W}_{\text{new}} = \left[\sum_{n=1}^N (\mathbf{x}_n - \bar{\mathbf{x}}) \mathbb{E}[\mathbf{z}_n]^T \right] \left[\sum_{n=1}^N \mathbb{E}[\mathbf{z}_n \mathbf{z}_n^T] \right]^{-1} \Rightarrow \mathbf{W}_{\text{new}} = \tilde{\mathbf{X}}^T \boldsymbol{\Omega}^T (\boldsymbol{\Omega} \boldsymbol{\Omega}^T)^{-1}$$

Re-estimation of the principal subspace minimizing the squared reconstruction errors in which the projections are fixed (see interpretation next and also [here](#))

EM Algorithm for PCA: Example, $D = 2, M = 1$

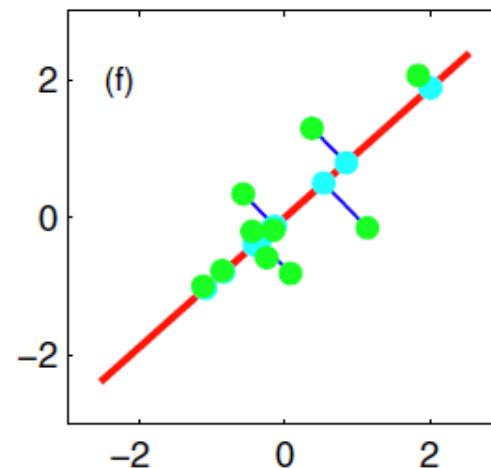
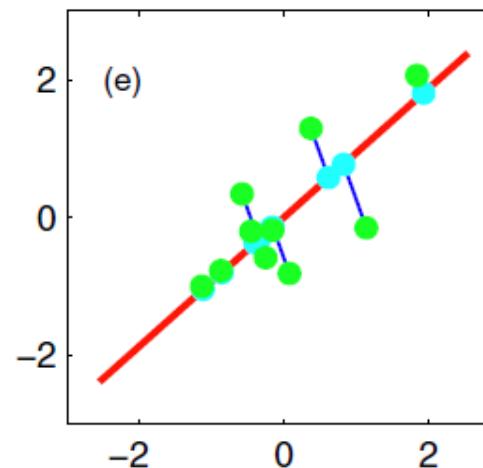
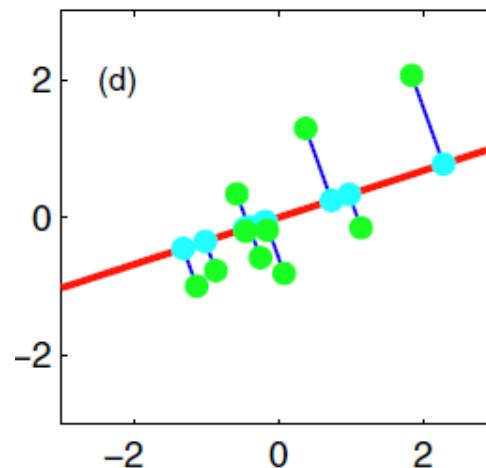
- Synthetic data illustrating the EM algorithm for PCA
 - (a) A data set X with the data points (green), together with the true principal components (eigenvectors scaled by the square roots of the eigenvalues).
 - (b) Initial configuration of the principal subspace defined by W (red) together with the projections of the latent points Z into the data space, given by ZW^T (cyan)
 - (c) After one M-Step, the latent space has been updated with Z held fixed.

[Matlab Implementation](#)



EM Algorithm for PCA: Example, $D = 2, M = 1$

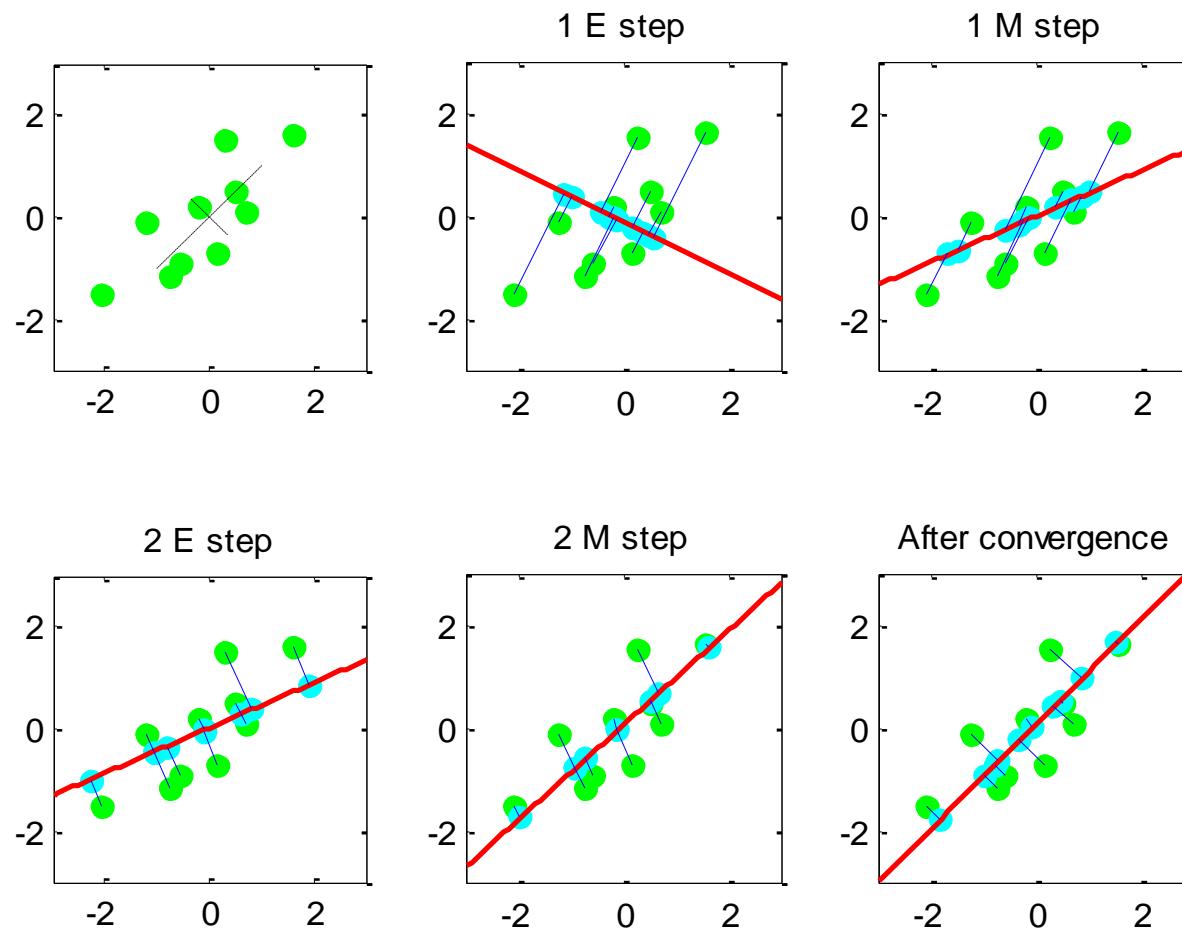
- Synthetic data illustrating the EM algorithm for PCA
 - (d) After the successive E-Step, the values of Z have been updated, giving orthogonal projections, with W held fixed.
 - (e) After the second M-Step.
 - (f) The converged solution.



$$\Omega = (\mathbf{W}_{old}^T \mathbf{W}_{old})^{-1} \mathbf{W}_{old}^T \tilde{\mathbf{X}}^T$$

$$\mathbf{W}_{new} = \tilde{\mathbf{X}}^T \Omega^T (\Omega \Omega^T)^{-1}$$

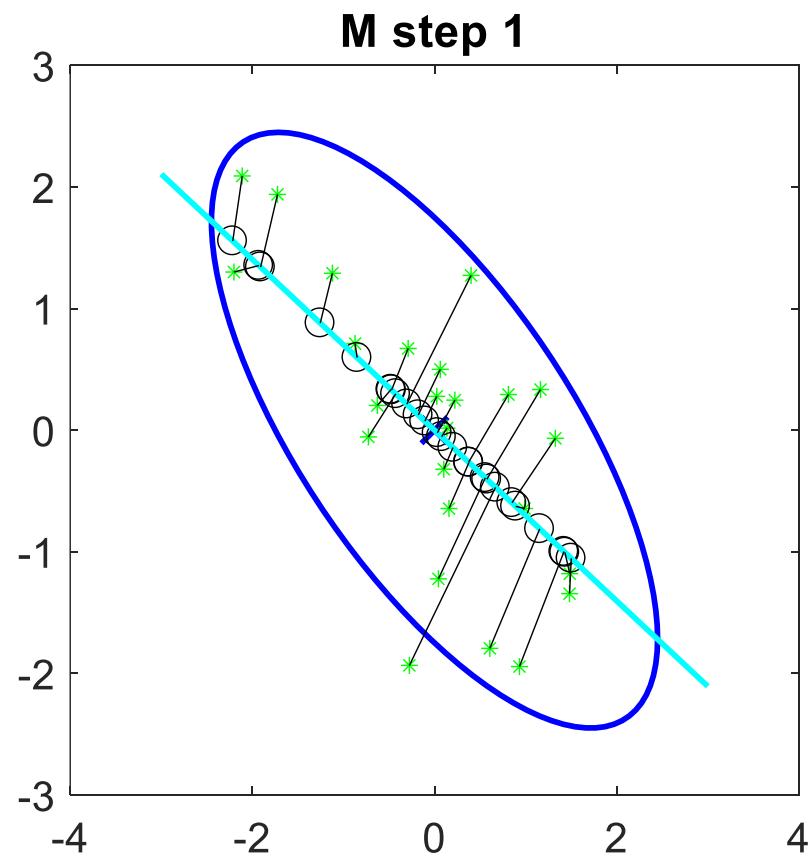
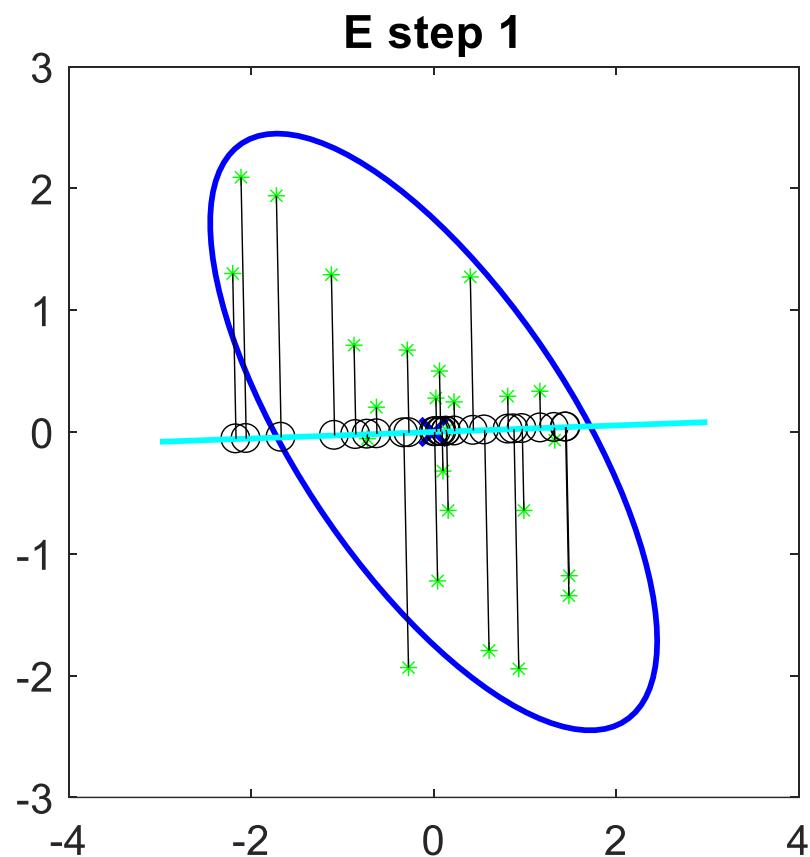
EM Algorithm for PCA: Example, $D = 2, M = 1$



[Matlab Implementation](#)



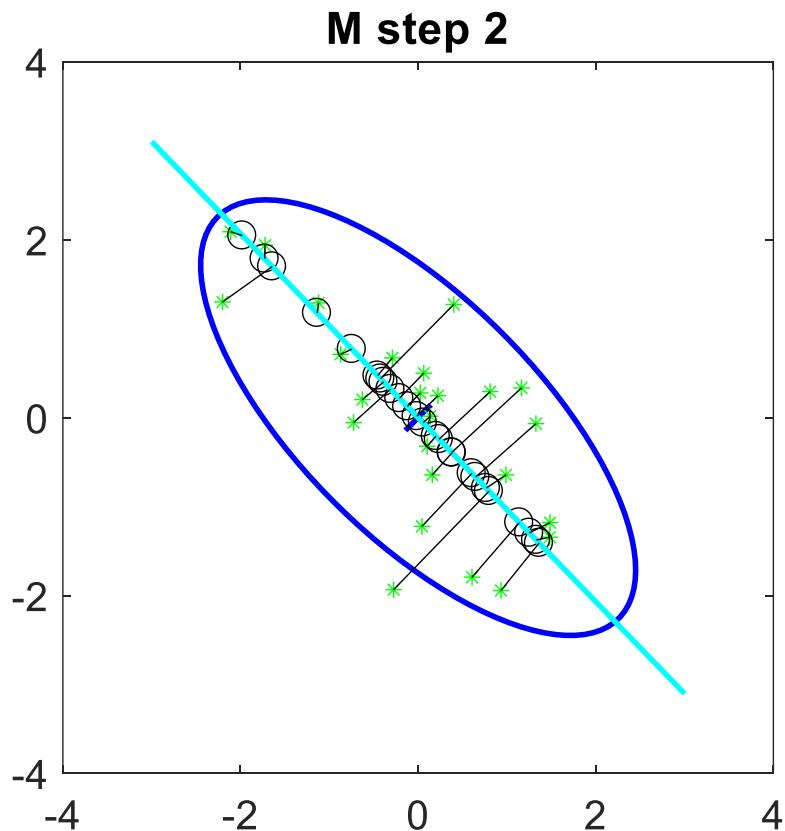
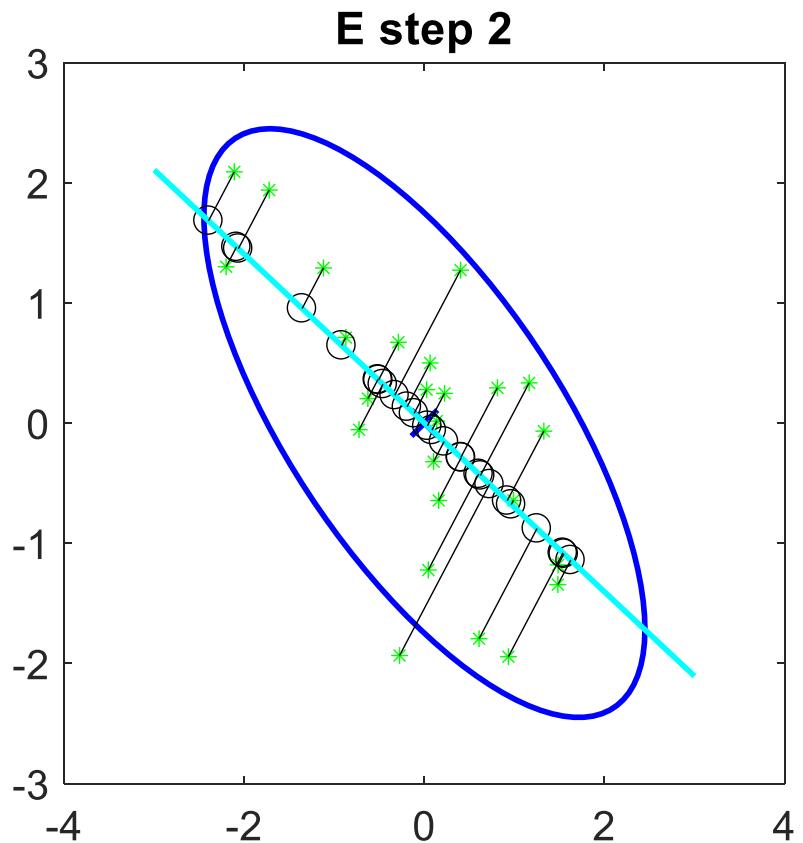
EM Algorithm for PCA: Example, $D = 2, M = 1$



Run [pcaEmStepByStep](#)
From [PMTK3](#)



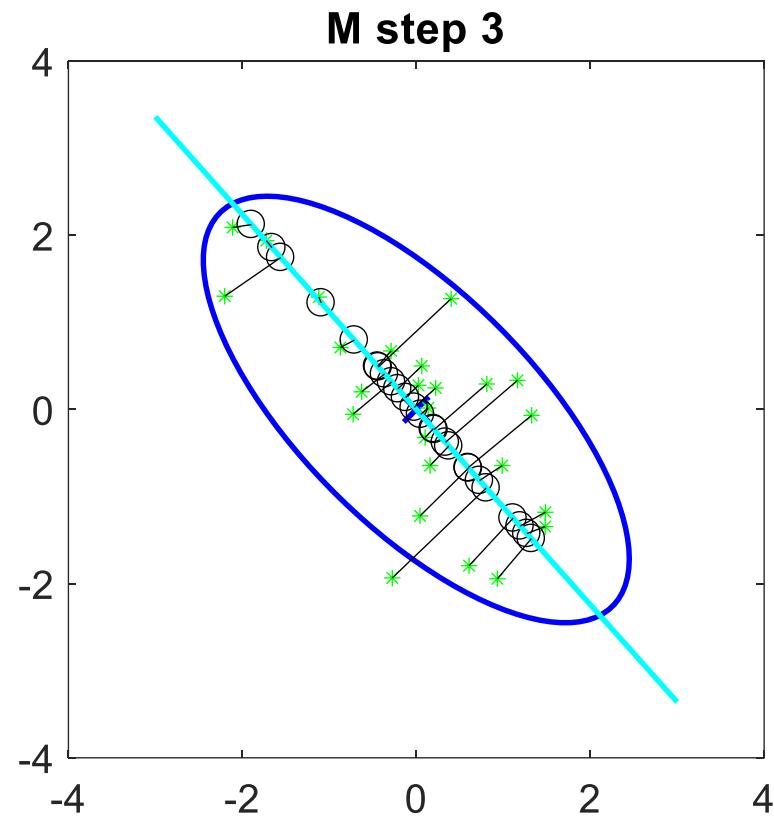
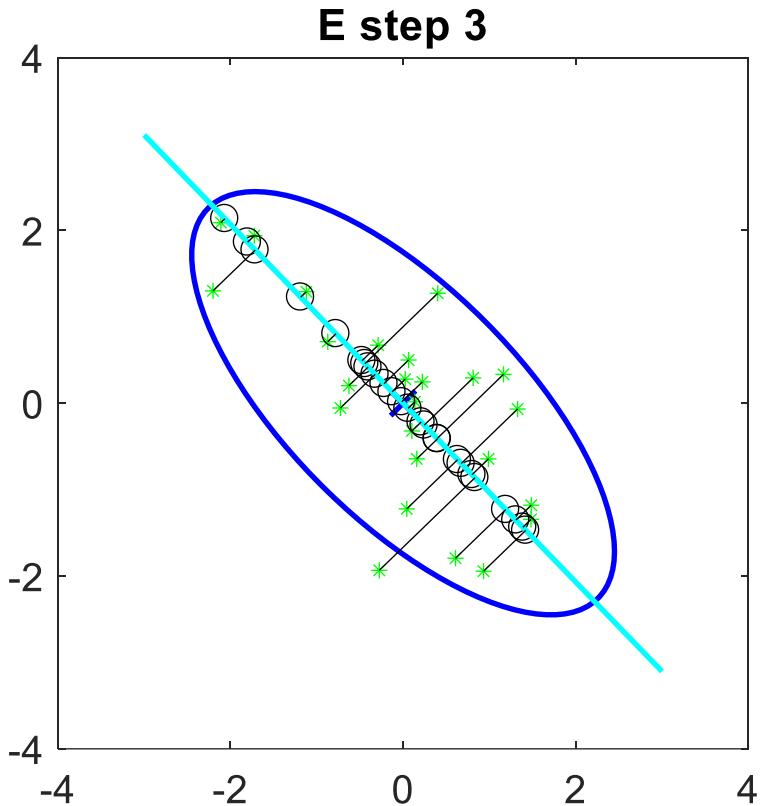
EM Algorithm for PCA: Example, $D = 2, M = 1$



Run [pcaEmStepByStep](#)
From [PMTK3](#)



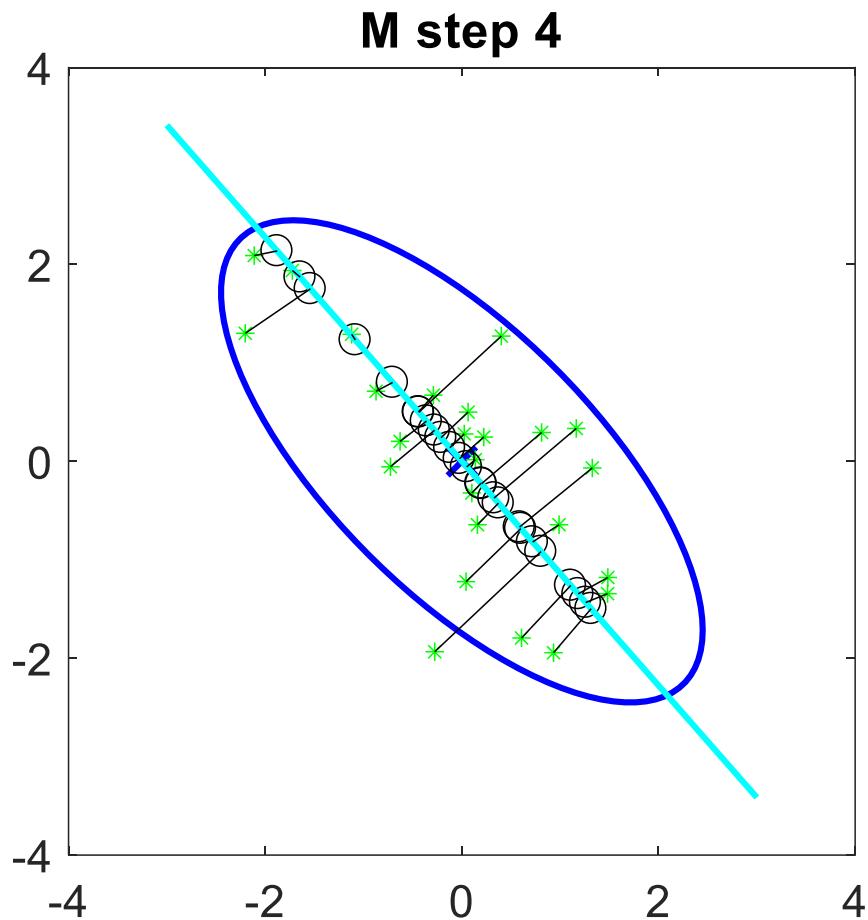
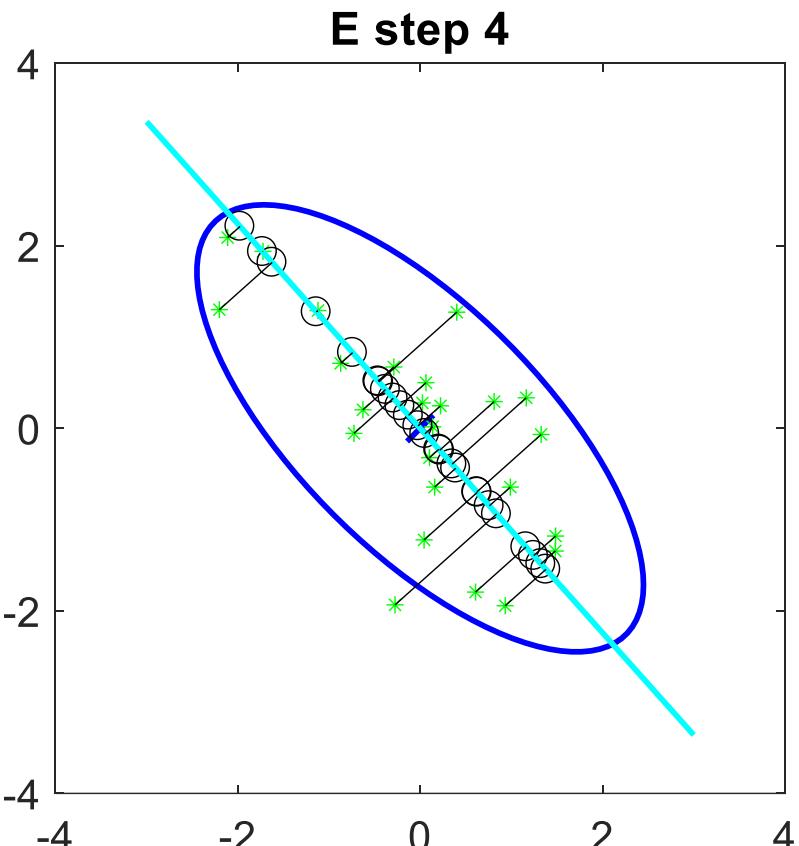
EM Algorithm for PCA: Example, $D = 2, M = 1$



Run [pcaEmStepByStep](#)
From [PMTK3](#)



EM Algorithm for PCA: Example, $D = 2, M = 1$



Run [pcaEmStepByStep](#)
From [PMTK3](#)



EM Algorithm and Standard PCA Revisited

- Let \mathbf{W} be $D \times M$ whose columns define a linear subspace of dimensionality M embedded within a data set of dimensionality D .
- Let $\boldsymbol{\mu}$ a D -dimensional vector. We approximate the data points $\mathbf{x}_n, n = 1, \dots, N$ using a linear mapping from a set of D -dimensional vectors \mathbf{z}_n as $\mathbf{Wz}_n + \boldsymbol{\mu}$.
- The sum of squares reconstruction error is:

$$J = \sum_{n=1}^N \|\mathbf{x}_n - \boldsymbol{\mu} - \mathbf{Wz}_n\|^2$$

- Minimizing wrt $\boldsymbol{\mu}$ gives:

$$0 = -\sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu} - \mathbf{Wz}_n) \Rightarrow \boldsymbol{\mu} = \bar{\mathbf{x}} - \mathbf{W}\bar{\mathbf{z}}$$

- This modifies the error as:

$$J = \sum_{n=1}^N \|(\mathbf{x}_n - \bar{\mathbf{x}}) - \mathbf{W}(\mathbf{z}_n - \bar{\mathbf{z}})\|^2$$



EM Algorithm and Standard PCA Revisited

- Let X be a $N \times D$ matrix with n^{th} row $(x_n - \bar{x})^T$.
- Similarly let Z be a $N \times M$ matrix with n^{th} row $(z_n - \bar{z})^T$.
- The cost function can be written as:

$$J = \sum_{n=1}^N \| (x_n - \bar{x}) - W(z_n - \bar{z}) \|^2 \Rightarrow J = \text{Tr} \left\{ (X - ZW^T)(X - ZW^T)^T \right\}$$

- Setting the derivative wrt Z equal to 0 gives (using $\frac{\partial}{\partial A} \text{Tr}(ABA^T) = A(B + B^T)$ and $\frac{\partial}{\partial A} \text{Tr}(AB) = B^T, \frac{\partial}{\partial A} \text{Tr}(A^T B) = B$) :

$$-2XW + 2ZW^TW = 0 \Rightarrow Z = XW(W^TW)^{-1}$$

This is similar to the PCA E-Step: $\Omega = (W_{old}^TW_{old})^{-1}W_{old}^T\tilde{X}^T$

- Similarly derivative wrt W equal to 0 gives (analogous to the PCA M-Step: $W_{new} = \tilde{X}^T\Omega^T(\Omega\Omega^T)^{-1}$)

$$-2X^TZ + 2WZ^TZ = 0 \Rightarrow W = X^TZ(Z^TZ)^{-1}$$



Nonlinear Latent Variable Models



Independent Component Analysis



Independent Component Analysis

- Setup:
 - Two people talking at the same time
 - Their voices recorded using two microphones
 - Objective: to reconstruct the 2 signals separately “blindly” (we are given only the mixed data). We haven’t observed
 - the original sources or the mixing coefficients
 - Under some assumptions (no time delay and echoes)
 - the signals received by the microphone are linear combinations of the voice amplitudes
 - the coefficient of this linear combination are constant
- Cardoso, J.-F. (1998). [Blind signal separation: statistical principles](#). *Proceedings of the IEEE* 9(10), 2009–2025.
- Stone, J. V. (2004). [Independent Component Analysis: A Tutorial Introduction](#). MIT Press.
- [A Tutorial on Independent Component Analysis](#) Jonathon Shlens
- A. Hyvärinen, E. Oja, [Independent component analysis: algorithms and applications](#), Neural Networks 13 (2000) 411–430
- A. Ng, ICA, [Lecture notes](#)
- [B. Poczos, Introduction to ICA](#), 2009



Independent Component Analysis

□ Model

$$x_1(t) = W_{11}z_1(t) + W_{12}z_2(t)$$
$$x_2(t) = W_{21}z_1(t) + W_{22}z_2(t)$$

□ Observations

$$\begin{pmatrix} x_1(1) \\ x_2(1) \end{pmatrix}, \begin{pmatrix} x_1(2) \\ x_2(2) \end{pmatrix}, \dots, \begin{pmatrix} x_1(t) \\ x_2(t) \end{pmatrix}$$

□ Objective: Estimate

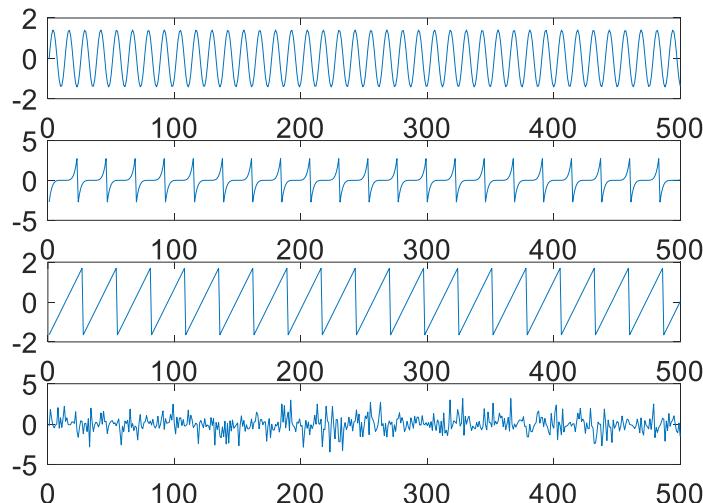
$$(z_1(1), z_2(1)), \dots (z_1(t), z_2(t))$$

and the coefficients W_{ij}



Independent Component Analysis

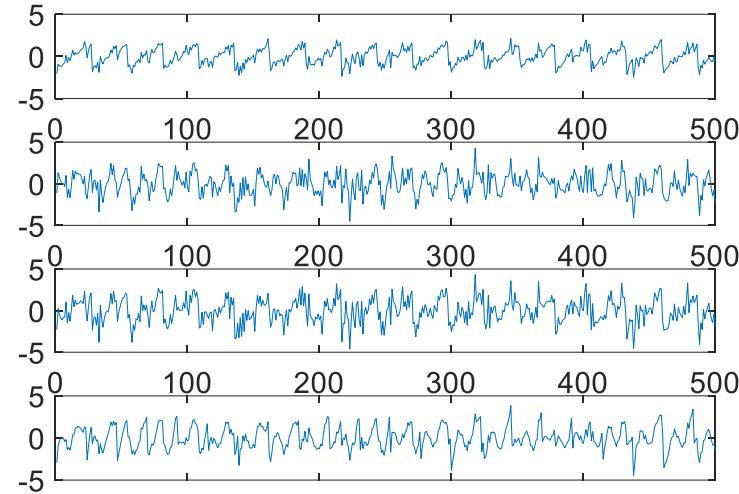
truth



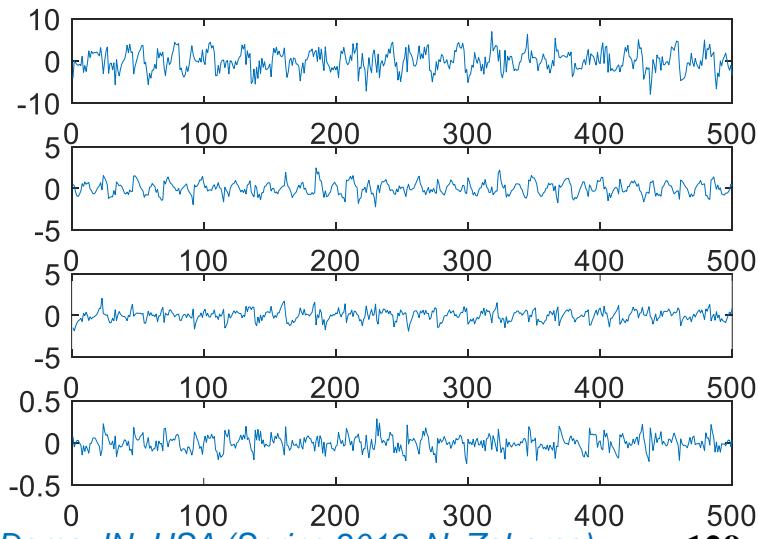
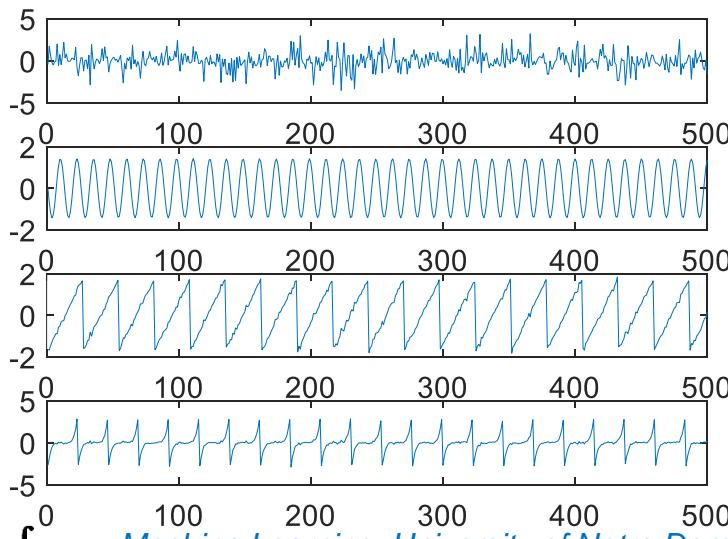
ICA estimate

[Run icaDemo](#)
From [PMTK3](#)

observed signals



PCA estimate



Independent Component Analysis

- Let $x_t \in \mathbb{R}^D$ be the observed signal at the sensors at time t , and $\mathbf{z}_t \in \mathbb{R}^L$ be the vector of source signals.

$$\mathbf{x}_t = \mathbf{W}\mathbf{z}_t + \boldsymbol{\epsilon}_t$$

- Here $\mathbf{W} \in \mathbb{R}^{D \times L}$, and $\boldsymbol{\epsilon}_t \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Psi})$. Often $\boldsymbol{\Psi} = \mathbf{0}$.
- We treat each time point as an independent observation (no temporal correlations).
- The goal is to infer the source signals, $p(\mathbf{z}_t | \mathbf{x}_t, \boldsymbol{\theta})$. \mathbf{W} is called the **mixing matrix**. If $L = D$ (number of sources = number of sensors), it will be a square matrix.
- The model up to now is identical to factor analysis (except we don't in general require orthogonality of \mathbf{W}).
- However, we will use a different prior for $p(\mathbf{z}_t)$. In PCA, we assume each source is independent, and has a Gaussian distribution

$$p(\mathbf{z}_t) = \prod_{j=1}^L \mathcal{N}(z_{tj} | 0, 1)$$



Independent Component Analysis

- In ICA, we let the sources be non-Gaussian distributions.

$$p(\mathbf{z}_t) = \prod_{j=1}^L p_j(z_{tj})$$

- We constrain the variance of the source distributions to be 1 (any other variance can be modelled by scaling the rows of \mathbf{W}).
- The resulting model is known as **independent component analysis** or **ICA**.
- Using Gaussian distribution as a source prior in ICA does not permit unique recovery of the sources (shown in the next slide).
- This is because the PCA likelihood is invariant to any orthogonal transformation of the sources \mathbf{z}_t and mixing matrix \mathbf{W} .
- PCA can recover the best linear subspace in which the signals lie, but cannot uniquely recover the signals themselves.



Difficulties with the ICA

- Let us return to the ICA problem with independent i.i.d. sources

$$\mathbf{x} = \mathbf{W}\mathbf{z}, \mathbf{z} = [z_1, \dots, z_L] \text{ are jointly independent}$$

- We can easily see that the sources can be recovered only up to permutation, sign and scale. Indeed note the following:

$$\mathbf{x} = (\mathbf{W}\mathbf{P}^{-1}\boldsymbol{\Lambda}^{-1})(\boldsymbol{\Lambda}\mathbf{P}\mathbf{z})$$

where \mathbf{P} is an arbitrary permutation matrix and $\boldsymbol{\Lambda}$ an arbitrary scaling matrix.

- In what follows we assume the following:

$$\mathbb{E}[\mathbf{z}] = \mathbf{0}, \text{ and } \mathbb{E}[\mathbf{z}\mathbf{z}^T] = \mathbf{I}$$

- The mixing matrix is not affected: $\mathbf{x} - \mathbb{E}[\mathbf{x}] = \mathbf{W}(\mathbf{z} - \mathbb{E}[\mathbf{z}])$

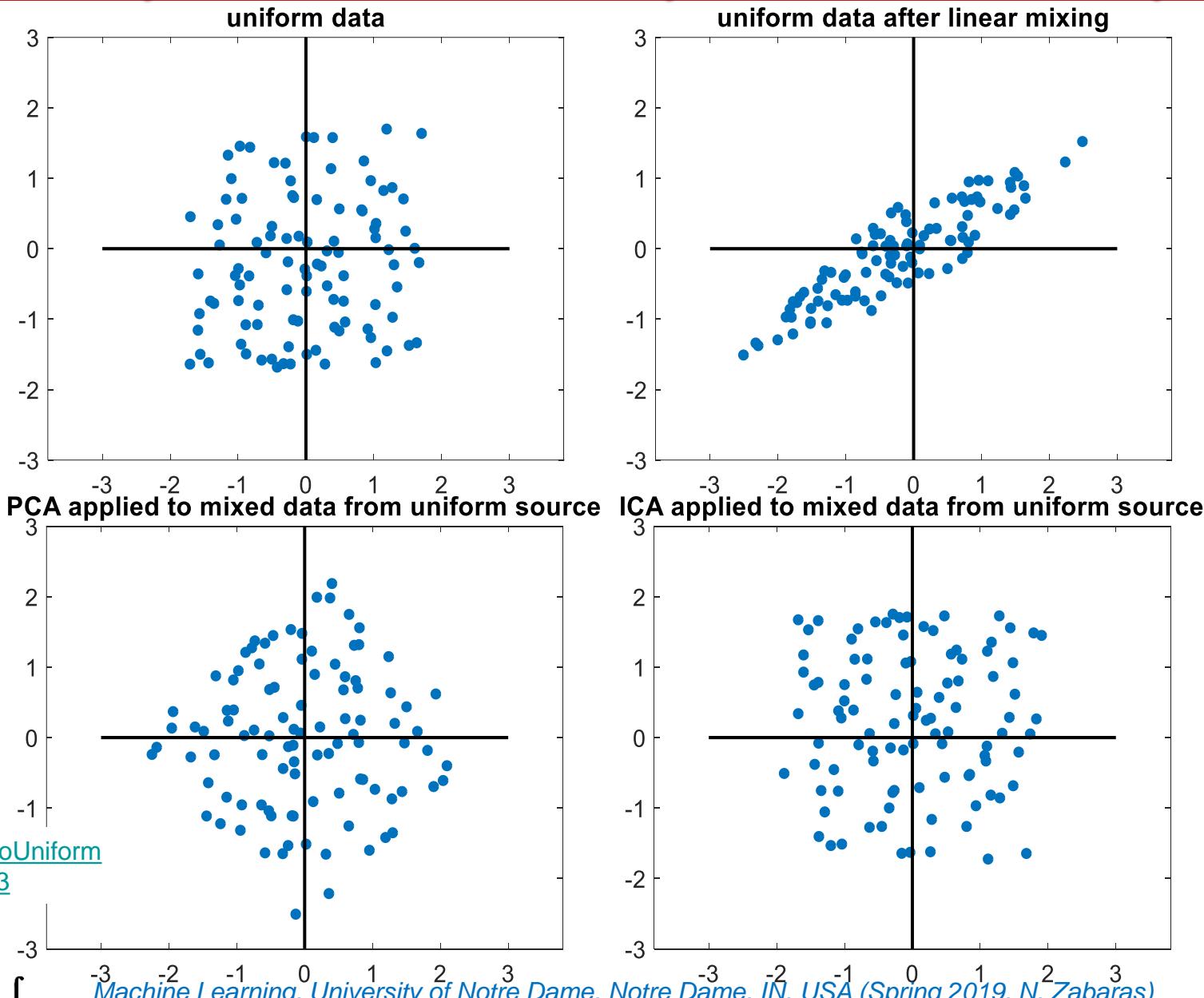


Independent Component Analysis

- Consider 2 independent sources with uniform PDFs with $W = \begin{pmatrix} 2 & 3 \\ 2 & 1 \end{pmatrix}$.
- Then we observe the data shown after this linear mixing.
- If we apply PCA followed by scaling, we get results that correspond to a whitening of the data.
- To uniquely recover the sources, we need to perform an additional rotation.
- However there is no information in the symmetric Gaussian posterior to tell us which angle to rotate by.
- PCA solves “half” of the problem, since it identifies the linear subspace; all that ICA has to do is then to identify the appropriate rotation.
- ICA is similar to varimax which seek good rotations of the latent factors to enhance interpretability.



Independent Component Analysis



Independent Component Analysis

- ICA can recover the source, up to a permutation of the indices and possible sign change.
- ICA requires that \mathbf{W} is square and hence invertible.
- In the non-square case (e.g., where we have more sources than sensors), we cannot uniquely recover the true signal, but we can compute the posterior $p(\mathbf{z}_t | \mathbf{x}_t, \widehat{\mathbf{W}})$.
- In both cases, we need to estimate \mathbf{W} as well as the source distributions p_j .

Review of Fundamental Concepts

Y_1, Y_2 are independent $\leftrightarrow p(Y_1, Y_2) = p(Y_1)p(Y_2)$

□ Let h_1, h_2 be arbitrary functions. The following holds:

Y_1, Y_2 are independent $\rightarrow \mathbb{E}[h_1(Y_1)h_2(Y_2)] = \mathbb{E}[h_1(Y_1)]\mathbb{E}[h_2(Y_2)]$

□ With the following definition of correlation

$$\text{corr}(Y_1, Y_2) = \frac{\mathbb{E}[(Y_1 - \mathbb{E}[Y_1])(Y_2 - \mathbb{E}[Y_2])]}{\sqrt{\text{var}[Y_1]}\sqrt{\text{var}[Y_2]}}$$

the following holds:

$$\text{corr}(Y_1, Y_2) = 0 \leftrightarrow \mathbb{E}[Y_1 Y_2] = \mathbb{E}[Y_1]\mathbb{E}[Y_2]$$

□ Let Y_1, Y_2 be jointly Gaussian. The following holds:

Y_1, Y_2 are independent $\leftrightarrow Y_1, Y_2$ are uncorrelated



Review of Fundamental Concepts

- Recall the definition of the **Shanon entropy**:

$$H(Y_1, Y_2, \dots, Y_L) = - \int p(y_1, y_2, \dots, y_L) \log(p(y_1, y_2, \dots, y_L)) dy$$

- The **mutual information** of a set of random variables is:

$$\begin{aligned} 0 \leq I(Y_1, Y_2, \dots, Y_L) &= \int p(y_1, y_2, \dots, y_L) \log \frac{p(y_1, y_2, \dots, y_L)}{p(y_1) \dots p(y_L)} dy \\ &= KL(p(y_1, y_2, \dots, y_L) || p(y_1) \dots p(y_L)) \\ &= \sum_{i=1}^L H(Y_i) - H(Y_1, Y_2, \dots, Y_L) \end{aligned}$$

Independent Component Analysis

- A possible approach is given in [Mackay'03](#) (chapter 34) that does not consider the temporal aspect of the problem and treats the successive samples as i.i.d.
- Consider generative model with
 - the two latent variables: unobserved speech signal amplitudes
 - the two observed signal values $x = [x_1 \ x_2]^T$ at the microphones
- Distribution of latent variables factorizes as $p(\mathbf{z}) = p(z_1)p(z_2)$
- No need to include noise: **observed variables = deterministic linear combinations of latent variables** as

$$x = \begin{bmatrix} W_{11} & W_{12} \\ W_{21} & W_{22} \end{bmatrix} \begin{pmatrix} z_1 \\ z_2 \end{pmatrix}$$

Blind Source Separation: ICA

- Given a set of observations
 - the likelihood function is a function of the coefficients W_{ij}
 - log likelihood maximized using gradient-based optimization

This is particular case of independent component analysis
- This requires that the latent variables have non Gaussian distributions
 - Probabilistic PCA: latent-space distribution = zero-mean isotropic Gaussian
 - No way to distinguish between two choices for the latent variables that differ by a rotation in the latent space
- PCA represents a rotation of the coordinate system in data space to diagonalize the covariance matrix. Zero correlation however does not imply independence!



Maximum Likelihood Estimation

- Consider estimating square \mathbf{W} for the noise-free ICA. We assume that the observations have been centered; hence we also assume \mathbf{z} is zero-mean. In addition, we assume the observations have been whitened (done with PCA).

$$\mathbb{E}[\mathbf{x}\mathbf{x}^T] = \mathbf{I}$$

- For the noise free case we have:

$$\text{cov}(\mathbf{x}) = \mathbb{E}[\mathbf{x}\mathbf{x}^T] = \mathbf{W}\mathbb{E}[\mathbf{z}\mathbf{z}^T]\mathbf{W}^T = \mathbf{W}\mathbf{W}^T$$

- Hence we see that \mathbf{W} must be orthogonal. This reduces the number of parameters we have to estimate from D^2 to $D(D - 1)/2$.
- Let $\mathbf{V} = \mathbf{W}^{-1}$ (**recognition weights**), as opposed to \mathbf{W} (**generative weights**)
- Since $\mathbf{x} = \mathbf{W}\mathbf{z}$, $p_{\mathbf{x}}(\mathbf{W}\mathbf{z}_t) = p_{\mathbf{z}}(\mathbf{z}_t)|\det(\mathbf{W}^{-1})| = p_{\mathbf{z}}(\mathbf{z}_t)|\det(\mathbf{V})|$. Thus for T iid samples (\mathbf{v}_j being the j -th row of \mathbf{V})

$$\frac{1}{T} \log p(\mathcal{D}|\mathbf{V}) = \log |\det(\mathbf{V})| + \frac{1}{T} \sum_{j=1}^L \sum_{t=1}^T \log p_j(\mathbf{v}_j^T \mathbf{x}_t)$$



Maximum Likelihood Estimation

$$\frac{1}{T} \log p(\mathcal{D}|V) = \log |\det(V)| + \frac{1}{T} \sum_{j=1}^L \sum_{t=1}^T \log p_j(v_j^T x_t)$$

- Since we are constraining V to be orthogonal, the 1st term is a constant. We can also replace the average over the data with an expectation.

$$NLL(V) = \sum_{j=1}^L \mathbb{E}[G_j(z_j)], z_j = v_j^T x, G_j(z) = -\log p_z(z)$$

- We want to minimize this subject to the constraint that the rows of V are orthogonal. We also want them to be unit norm, since this ensures that the variance of the factors is unity (with whitened data, $\mathbb{E}[(v_j^T x)(x^T v_j)] = \|v_j\|^2$ which is necessary to fix the scale of the weights. In otherwords, V should be an orthonormal matrix).
- Possible approaches include:
 - ✓ gradient descent algorithm to fit this model (rather slow);
 - ✓ natural gradient (faster) ([MacKay 2003](#), ch 34);
 - ✓ Approximate Newton method; EM.

Fast ICA Algorithm

- Define the following:

$$\begin{aligned}f(\boldsymbol{\nu}) &= \mathbb{E}[G(\boldsymbol{\nu}^T \mathbf{x})] + \lambda(1 - \boldsymbol{\nu}^T \boldsymbol{\nu}) \\ \nabla f(\boldsymbol{\nu}) &= \mathbb{E}[\mathbf{x}g(\boldsymbol{\nu}^T \mathbf{x})] - \beta\boldsymbol{\nu}, \beta = 2\lambda \\ \mathbf{H}(\boldsymbol{\nu}) &= \mathbb{E}[\mathbf{x}\mathbf{x}^T g'(\boldsymbol{\nu}^T \mathbf{x})] - \beta\mathbf{I}\end{aligned}$$

- Consider the following approximation:

$$\mathbb{E}[\mathbf{x}\mathbf{x}^T g'(\boldsymbol{\nu}^T \mathbf{x})] \approx \mathbb{E}[\mathbf{x}\mathbf{x}^T]\mathbb{E}[g'(\boldsymbol{\nu}^T \mathbf{x})] = \mathbb{E}[g'(\boldsymbol{\nu}^T \mathbf{x})]$$

- The following Newton update can now be introduced:

$$\boldsymbol{\nu}^* = \boldsymbol{\nu} - \frac{\mathbb{E}[\mathbf{x}g(\boldsymbol{\nu}^T \mathbf{x})] - \beta\boldsymbol{\nu}}{\mathbb{E}[g'(\boldsymbol{\nu}^T \mathbf{x})] - \beta} \text{ or } \boldsymbol{\nu}^* \triangleq \mathbb{E}[\mathbf{x}g(\boldsymbol{\nu}^T \mathbf{x})] - \mathbb{E}[g'(\boldsymbol{\nu}^T \mathbf{x})]\boldsymbol{\nu}$$

- The expectations can be replaced by Monte Carlo estimates from the training set (efficient online learning algorithm). After performing this update, one should project back onto the constraint surface using

$$\boldsymbol{\nu}^{new} = \frac{\boldsymbol{\nu}^*}{\|\boldsymbol{\nu}^*\|}$$



Fast ICA Algorithm

- Define the following:

$$\boldsymbol{v}^{new} = \frac{\boldsymbol{v}^*}{\|\boldsymbol{v}^*\|}$$

- One iterates this algorithm until convergence.

- ✓ Due to the sign ambiguity of \boldsymbol{v} , the values of \boldsymbol{v} may not converge, but the direction defined by this vector should converge, so
 - ✓ can assess convergence by monitoring $|\boldsymbol{v}^T \boldsymbol{v}^{new}|$, which should approach 1

- Since the objective is not convex, there are multiple local optima.



Fast ICA Algorithm

- We can learn multiple different weight vectors or **features**.
 - We can either learn the features sequentially and then project out the part of v_j that lies in the subspace defined by earlier features, or
 - We can learn them in parallel, and **orthogonalize V** in parallel.
 - This latter approach is usually preferred, since, unlike PCA, the features are not ordered in any way (e.g. the 1st feature is not more important than the 2nd and hence it is better to treat them symmetrically).



Modeling the Source Densities

- So far, we have assumed that $G(z) = -\log p(z)$ is known.
- What kinds of models might be reasonable as signal priors?
- We know that using Gaussians (which correspond to quadratic functions G) won't work.
- So we want some kind of non-Gaussian distribution.
- In general, there are several kinds of non-Gaussian distributions, such as the following:
 - Super-Gaussian distributions: $\text{kurt}(z) > 0$, e.g., Laplace distribution.
 - Sub-Gaussian distribution: $\text{kurt}(z) < 0$. e.g., Uniform distribution.
 - Skewed distribution: Another way to “be non-Gaussian” is to be asymmetric (one possible measure of asymmetry is skewness). E.g., Gamma distribution (right skewed).

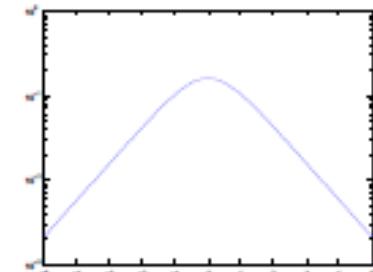


Independent Component Analysis

- Common choice for latent-variable distribution:

$$p(z_j) = \frac{1}{\pi \cosh(z_j)} = \frac{2}{\pi(e^{z_j} + e^{-z_j})} \text{ (heavy tails)}$$

- Original ICA model ([Bell and Sejnowski, 1995](#)) based on optimization of an objective function defined by information maximization.



- Bell, A. J. and T. J. Sejnowski (1995). [An information maximization approach to blind separation and blind deconvolution. *Neural Computation* 7\(6\), 1129–1159.](#)

Source Densities

- SuperGaussian (big spike at the mean, heavy tails, thin kurtosis): The Laplace distribution is an example. With μ_k denoting the k^{th} moment of the distribution, we define:

$$\text{kurt}(z) = \frac{\mu_4}{\sigma^2} - 3 > 0$$

We subtract 3 from the standard definition to make the kurtosis of a Gaussian distribution equal to zero.

- SubGaussian (broad kurtosis): $\text{kurt}(z) < 0$
- Skewed distributions: The skewness is a measure of the degree of the asymmetry of the distribution:

$$\text{skew}(z) = \frac{\mu_3}{\sigma^3}$$

- Pham, D.-T. and P. Garrat (1997). Blind separation of mixture of independent sources through a quasimaximum likelihood approach. *IEEE Trans. on Signal Processing* 45(7), 1712–1725.



Source Densities

- The empirical distribution of many natural signals (images, speech, when passed through certain linear filters) tend to be super-Gaussian.
- This holds for the kind of linear filters found in the brain (e.g. visual cortex), and in signal processing (wavelet transforms).
- One obvious choice for modeling natural signals with ICA is therefore the Laplace distribution. For mean zero and variance 1, this has a log pdf given by

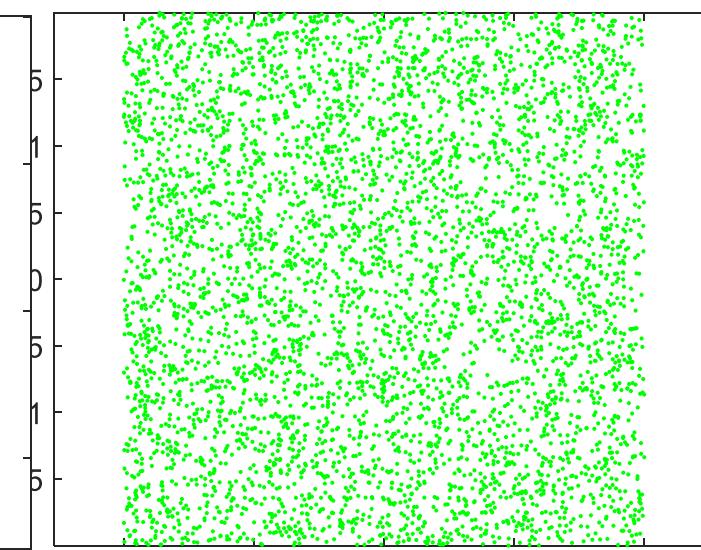
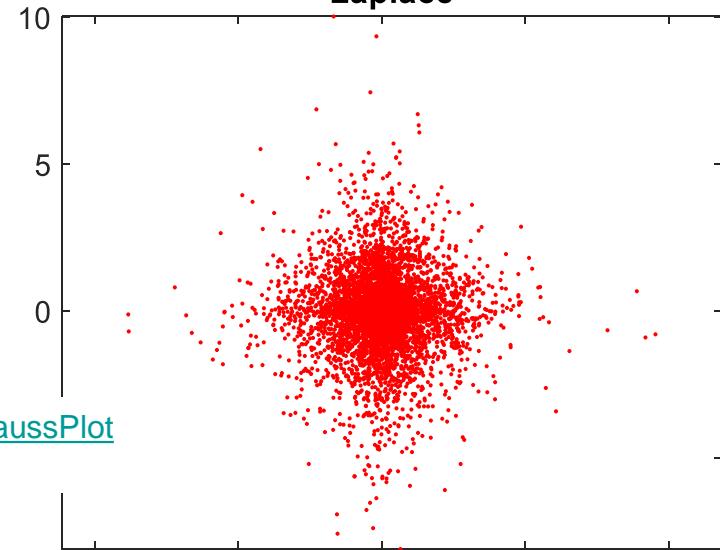
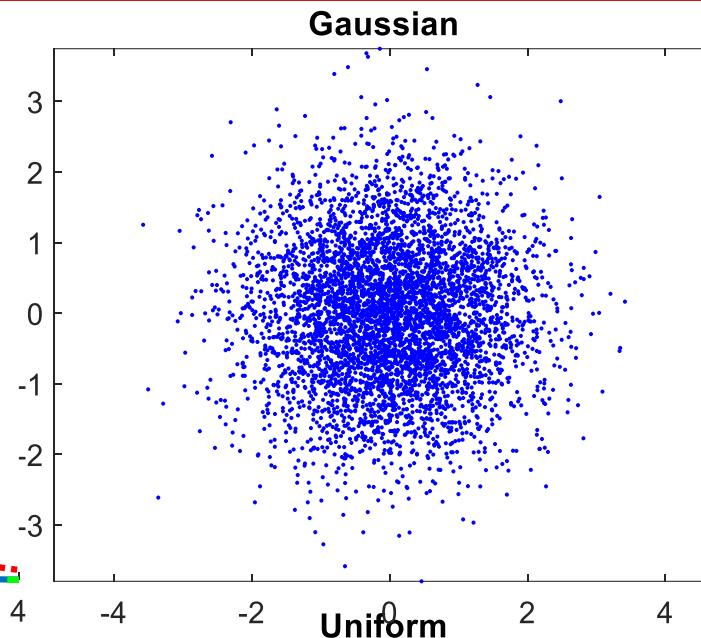
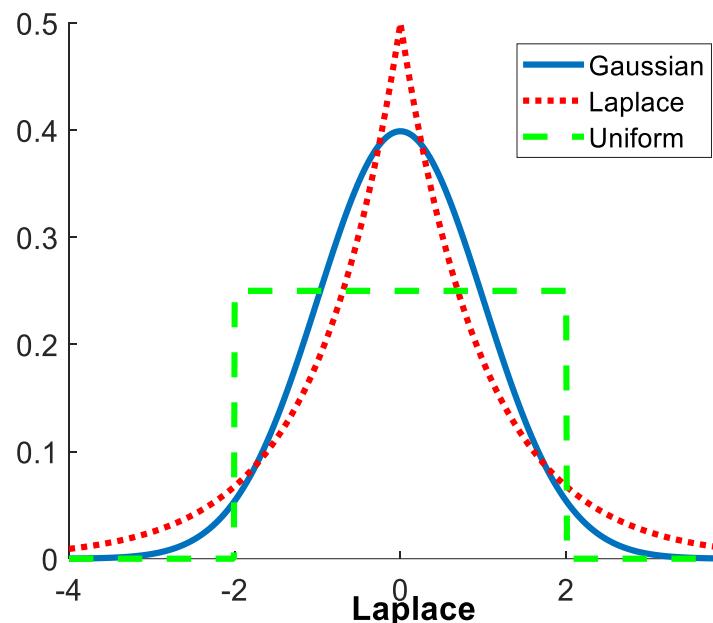
$$\log p(z) = -\sqrt{2}|z| - \log(\sqrt{2})$$

- The logistic distribution is a smoother super-Gaussian distribution. For mean 0 and variance 1 (so $\mu = 0, s = \frac{\sqrt{3}}{\pi}$):

$$\log p(z) = -2 \operatorname{logcosh} \left(\frac{\pi}{2\sqrt{3}} z \right) - \log \left(\frac{4\sqrt{3}}{\pi} \right)$$

- As the exact type of source distribution is not that important, we often use: $G(z) = -\log p(z) = |z|$ or $G(z) = \operatorname{logcosh}(z)$.

Source Densities



Run subSuperGaussPlot
From PMTK3

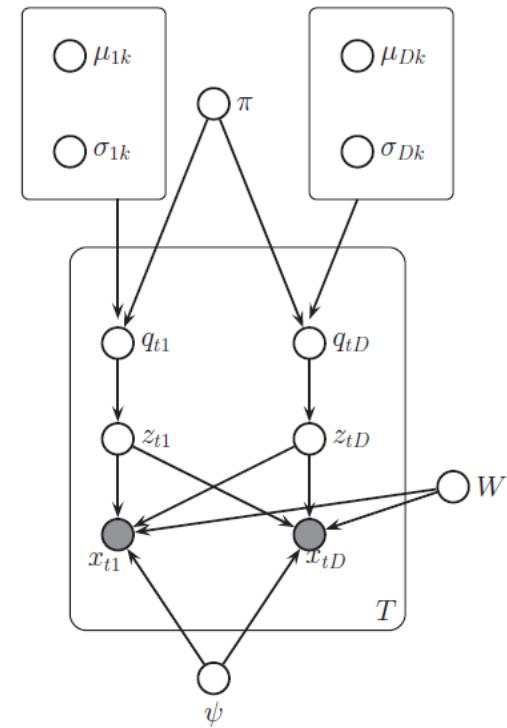


Using EM

- An alternative to assuming a form for $p(\mathbf{z})$, is to use a mixture of (uni-variate) Gaussians:

$$\begin{aligned} p(q_j = k) &= \pi_k \\ p(z_j | q_j = k) &= \mathcal{N}(\mu_j, \sigma_{jk}^2) \\ p(\mathbf{x}|\mathbf{z}) &= \mathcal{N}(\mathbf{Wz}, \boldsymbol{\Psi}) \end{aligned}$$

- We can compute $\mathbb{E} [\mathbf{z}_t | \mathbf{x}_t, \theta]$ by summing over all K^L combinations of the q_t variables, where K is the number of mixture components per source (one can also use a mean field approximation)
- We can then estimate all the source distributions in parallel by fitting a standard GMM to $\mathbb{E}[\mathbf{z}_t]$. When the source GMMs are known, we can compute the marginals $p_j(z_j)$ very easily, using $p_j(z_j) = \sum_{k=1}^K \pi_{jk} \mathcal{N}(\mu_{jk}, \sigma_{jk}^2)$
- Given the π_j 's, we can then use an ICA algorithm to estimate \mathbf{W} .
- These steps should be interleaved.



- Attias, H. (1999a). [Independent factor analysis](#). *Neural Computation* 11(4), 803–851.
- Moulines, E., J.-F. Cardoso, and E. Gassiat (1997). [Maximum likelihood for blind separation and deconvolution of noisy signals using mixture models](#). In *Proc. IEEE Int. Conf. on Acoustics, Speech and Signal Processing (ICASSP'97)*, Munich Germany, pp. 3617–3620.

Independent Component Analysis

- Independent factor analysis ([Attias, 1999](#)) considers a model in which the number of latent and observed variables can differ, the observed variables are noisy, and the individual latent variables have flexible distributions modeled by mixtures of Gaussians.
 - The log likelihood for this model is maximized using EM, and the reconstruction of the latent variables is approximated using a variational approach.
- References to other works are given next

- Attias, H. (1999a). [Independent factor analysis](#). *Neural Computation* **11**(4), 803–851.



Independent Component Analysis

□ Some key works include:

- Hyvarinen, A. and E. Oja (1997). [A fast fixed-point algorithm for independent component analysis](#). *Neural Computation* **9**(7), 1483–1492.
- Attias, H. (1999a). [Independent factor analysis](#). *Neural Computation* **11**(4), 803–851.
- Jutten, C. and J. Herault (1991). [Blind separation of sources, 1: An adaptive algorithm based on neuromimetic architecture](#). *Signal Processing* **24**(1), 1–10.
- Comon, P., C. Jutten, and J. Herault (1991). [Blind source separation, 2: problems statement](#). *Signal Processing* **24**(1), 11–20.
- Amari, S., A. Cichocki, and H. H. Yang (1996). [A new learning algorithm for blind signal separation](#). In D. S. Touretzky, M. C. Mozer, and M. E. Hasselmo (Eds.), [Advances in Neural Information Processing Systems](#), Volume 8, pp. 757–763. MIT Press.
- Pearlmutter, B. A. and L. C. Parra (1997). [Maximum likelihood source separation: a context-sensitive generalization of ICA](#). In M. C. Mozer, M. I. Jordan, and T. Petsche (Eds.), [Advances in Neural Information Processing Systems, Volume 9](#), pp. 613–619. MIT Press.
- Hinton, G. E., M. Welling, Y. W. Teh, and S. Osindero (2001). [A new view of ICA](#). In [Proceedings of the International Conference on Independent Component Analysis and Blind Signal Separation](#), Volume 3.
- Miskin, J. W. and D. J. C. MacKay (2001). [Ensemble learning for blind source separation](#). In S. J. Roberts and R. M. Everson (Eds.), [Independent Component Analysis: Principles and Practice](#). Cambridge University Press.
- Hojen-Sorensen, P. A., O. Winther, and L. K. Hansen (2002). [Mean field approaches to independent component analysis](#). *Neural Computation* **14**(4), 889–918.
- Choudrey, R. A. and S. J. Roberts (2003). [Variational mixture of Bayesian independent component analyzers](#). *Neural Computation* **15**(1), 213–252.
- Chan, K., T. Lee, and T. J. Sejnowski (2003). [Variational Bayesian learning of ICA with missing data](#). *Neural Computation* **15**(8), 1991–2011.
- Stone, J. V. (2004). [Independent Component Analysis: A Tutorial Introduction](#). MIT Press.



Other Estimation Principles

- It is quite common to estimate the parameters of ICA models using methods that seem different to maximum likelihood.
- We will review some of these methods below, because they give additional insight into ICA.
- However, we will also see that these methods in fact are equivalent to maximum likelihood after all.

Maximizing non-Gaussianity

- Can pose the ICA problem as follows: Find \mathbf{V} such that the distribution $\mathbf{z} = \mathbf{V}\mathbf{x}$ is as far from Gaussian as possible.
- One measure of non-Gaussianity is kurtosis, but this can be sensitive to outliers. Another measure is the **negentropy**, defined as

$$\text{negentropy}(\mathbf{z}) = \mathbb{H}(\mathcal{N}(\mu, \sigma^2)) - \mathbb{H}(\mathbf{z}), \mu = \mathbb{E}[\mathbf{z}] \text{ and } \sigma^2 = \text{var}[\mathbf{z}].$$

- Since the Gaussian is the max entropy distribution, this measure is always non-negative and becomes large for distributions that are highly non-Gaussian. Using $\mathbf{z} = \mathbf{V}\mathbf{x}$, \mathbf{V} orthogonal, whitten data (so $\text{cov}(\mathbf{z}) = \mathbf{I}$, independent of \mathbf{V}), we can define our objective as maximizing

$$J(\mathbf{V}) = \sum_j \text{negentropy}(z_j) = \sum_j \{\mathbb{H}(\mathcal{N}(\mu_j, \sigma_j^2)) - \mathbb{H}(z_j)\} = \sum_j -\mathbb{H}(z_j) + \text{const}$$

- Considering $\mathbb{H}(z_j) = -\mathbb{E}[\log p(z_j)]$, this gives the same form as maximizing the log-likelihood.

- Hyvärinen, A. and E. Oja (2000). [Independent component analysis: algorithms and applications](#). *Neural Networks* 13, 411–430.



Minimizing Mutual Information

- Since we are trying to find independent components, we can minimize the following:

$$I(\mathbf{z}) = KL\left(p(\mathbf{z}) \parallel \prod_j p(z_j)\right) = \sum_j -\mathbb{H}(z_j) - \mathbb{H}(\mathbf{z}) = -\sum_j \mathbb{H}(z_j) - \mathbb{H}(\mathbf{V}\mathbf{x})$$

- For \mathbf{V} orthogonal, we can drop the last term, since then $\mathbb{H}(\mathbf{V}\mathbf{x}) = \mathbb{H}(\mathbf{x})$ (multiplying by \mathbf{V} does not change the shape of the distribution), and $H(\mathbf{x})$ is a constant which is solely determined by the empirical distribution.
- Hence, we have $I(\mathbf{z}) = -\sum_j \mathbb{H}(z_j) + const.$
- Minimizing this is equivalent to maximizing the negentropy, which is equivalent to maximum likelihood.

Maximizing Mutual Information

- Imagine a neural network where x is the input and $y_j = \phi(v_j^T x) + \varepsilon$ is the noisy output, ϕ is some nonlinear scalar function, and $\varepsilon \sim \mathcal{N}(0, 1)$.
- We want to maximize the mutual information between y (the internal neural representation) and x (the observed input signal). We have $I(x; y) = H(y) - H(y|x)$, where the latter term is constant if we assume the noise has constant variance. One can show:

$$H(y) = \sum_{j=1}^L \mathbb{E}[\log \phi'(v_j^T x)] + \log |\det(V)|$$

- We can drop the last term if V is orthogonal. If we define $\phi(z)$ to be a cdf, then $\phi'(z)$ is its pdf, and the above expression is equivalent to the log likelihood.
- If we use a logistic nonlinearity, $\phi(z) = \text{sigm}(z)$, the corresponding pdf is the logistic distribution, and $\log \phi'(z) = \log \cosh(z)$. Thus **we see that infomax is equivalent to maximum likelihood**.

- Bell, A. J. and T. J. Sejnowski (1995). [An information maximisation approach to blind separation and blind deconvolution](#). *Neural Computation* 7(6), 1129–1159.

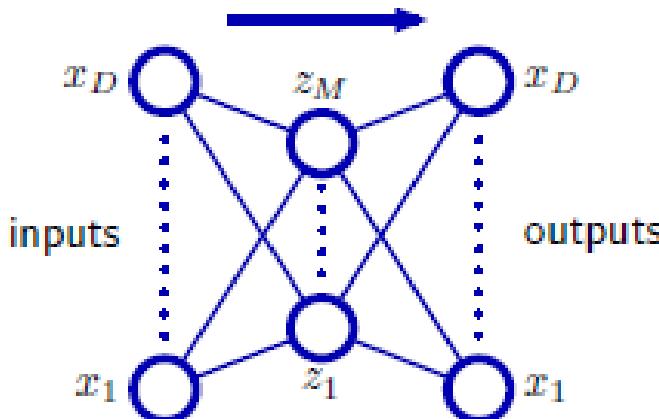


Neural Networks for Dimensionality Reduction



Autoassociative Neural Networks

- Neural networks for predicting outputs given inputs, can also be used for dimensionality reduction.

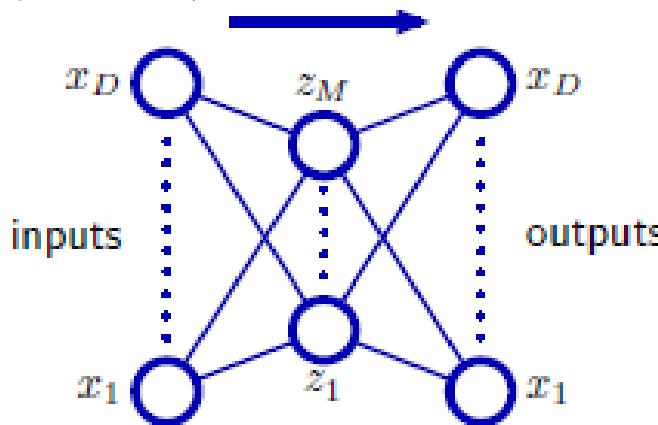


- We need a network having the same number of outputs as inputs, and optimize the weights to minimize the reconstruction error between inputs and outputs with respect to the training data set:

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \|y(x_n, \mathbf{w}) - x_n\|^2$$

Autoassociative Neural Networks

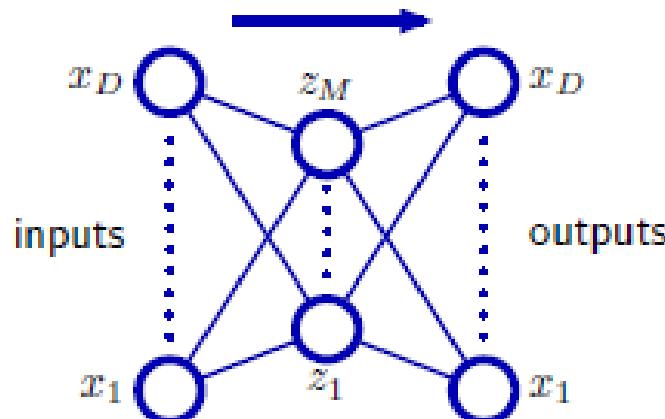
- Consider a multilayer perceptron: D inputs, D outputs and M hidden units, with $M < D$.



- The targets used to train the network are the input vectors themselves i.e. the network is attempting to map each input vector onto itself (autoassociative mapping).
- Since $M < D$, a perfect reconstruction of all input vectors is not in general possible.
- Determine w by minimizing the mismatch between the input vectors and their reconstructions.

Autoassociative Neural Networks

- Use the sum-of-square errors



$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \|y(x_n, \mathbf{w}) - x_n\|^2$$

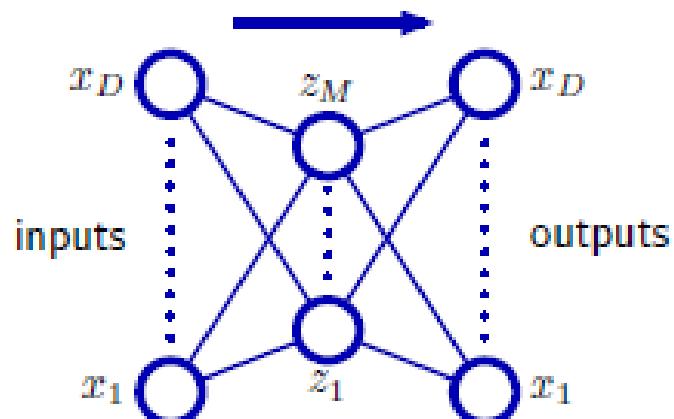
- For hidden units with linear activation functions,

- the error function has a unique global minimum, and
- at this minimum the network performs a projection onto the M -dimensional subspace which is spanned by the first M principal components of the data.

- [Bourlard, H. and Y. Kamp \(1988\). Auto-association by multilayer perceptrons and singular value decomposition. *Biological Cybernetics* 59](#), 291– 294.
- [Baldi, P. and K. Hornik \(1989\). Neural networks and principal component analysis: learning from examples without local minima. *Neural Networks* 2\(1\), 53–58.](#)

Autoassociative Neural Networks

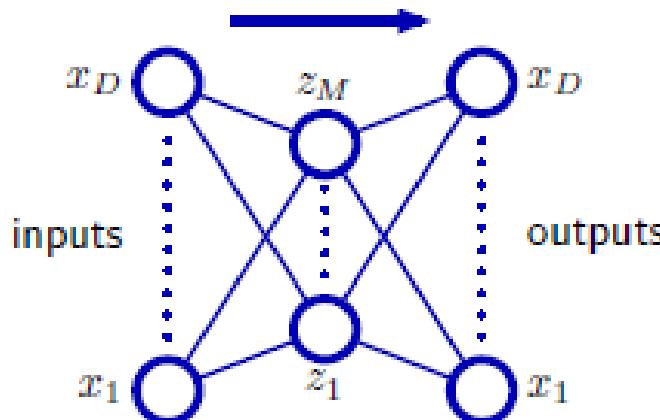
- Thus, the vectors of weights which lead into the hidden units form a basis set which spans the principal subspace. These vectors need not be orthogonal or normalized.
- Both PCA and the neural network use linear dimensionality reduction and are minimizing the same sum-of-squares error function → they both perform projection onto the M –dimensional subspace spanned by the first M –principal components of the data:



$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \|y(x_n, \mathbf{w}) - x_n\|^2$$

Autoassociative Neural Networks

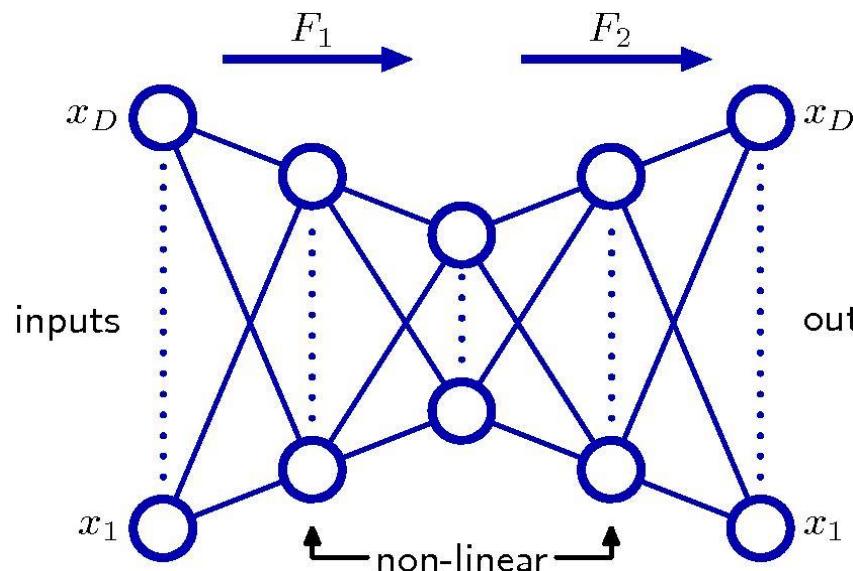
- Surprisingly, even with nonlinear hidden units in the two layers network below, the minimum error solution is again given by the projection onto the principal component subspace ([Bourlard and Kamp, 1988](#)).



- There is therefore no advantage in using two layer neural networks to perform dimensionality reduction.
 - Standard PCA gives the correct solution in finite time, and generates ordered eigenvalues with orthonormal eigenvectors.
 - [Bourlard, H. and Y. Kamp \(1988\). Auto-association by multilayer perceptrons and singular value decomposition. Biological Cybernetics 59](#), 291– 294.

Autoassociative Neural Networks

- ❑ Use more hidden layers (4 here) with some non-linear activation functions.

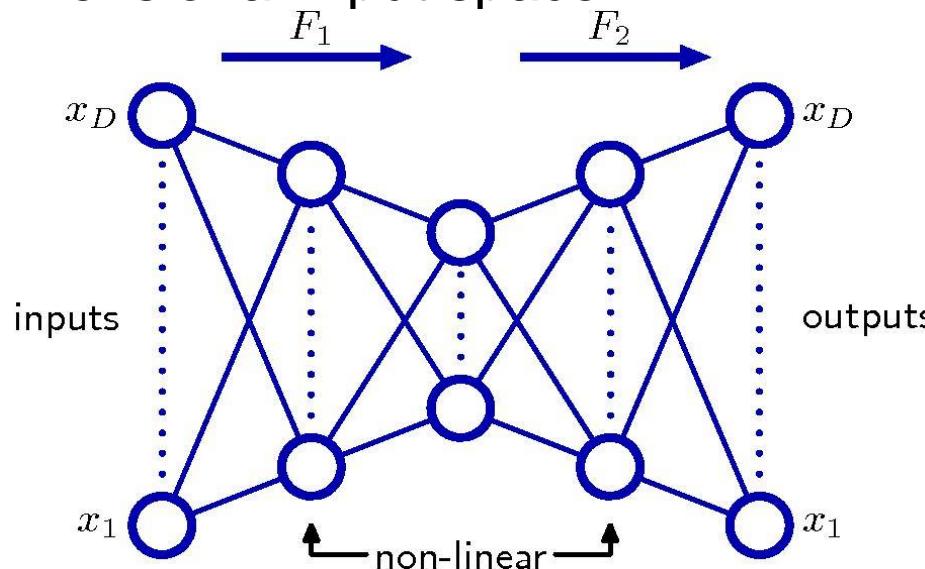


$$\text{outputs } E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \|y(x_n, \mathbf{w}) - x_n\|^2$$

- ❑ The output units are linear, and the M units in the 2nd hidden layer can also be linear. However, the 1st and 3rd hidden layers have sigmoidal nonlinear activation functions.
- ❑ Training involves **minimization of the error function with risk of suboptimality**.

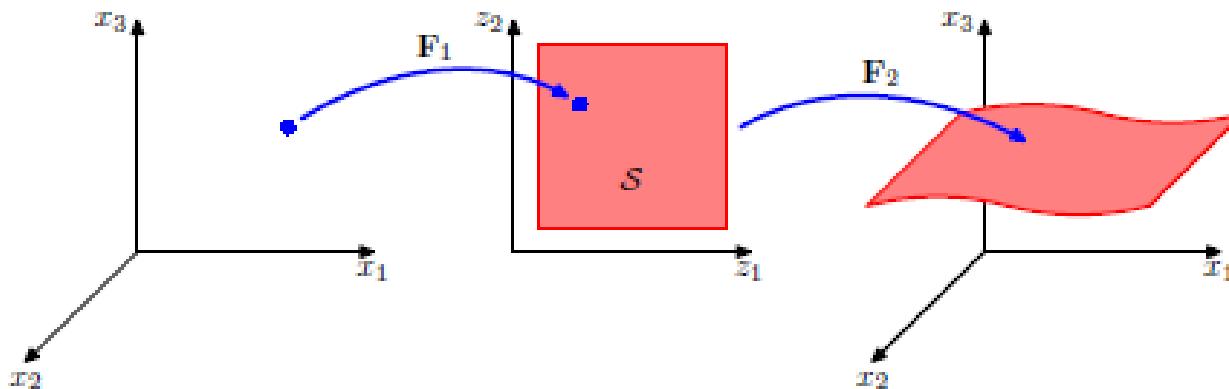
Autoassociative Neural Networks

- The map F_1 projects the original D dimensional data onto an M –dimensional subspace S defined by the activations of the units in the 2nd hidden layer.
- Because of the presence of the first hidden layer of nonlinear units, this mapping is not linear.
- Similarly, the 2nd half of the network defines an arbitrary functional mapping from the M –dimensional space back into the original D –dimensional input space.



Nonlinear PCA with Autoassociative Neural Networks

- This has a simple geometrical interpretation. Consider $D = 3$ and $M = 2$.



- F_2 maps the M –dimensional S into a D –dimensional space and therefore defines the way in which the space S is embedded within the original x –space.
- Since F_2 can be nonlinear, the embedding of S can be nonplanar.
- F_1 defines a projection of points from the original D –dimensional space into the M –dimensional subspace S .

Modeling Nonlinear Manifolds



Modeling NonLinear Manifolds

- Data often lie in a manifold of lower-dimensionality than the observed data space
- Capturing this may improve the density modeling
- Possible approach: non-linear manifold modeled by piece-wise linear approximation, e.g.,
 - k-means + PCA for each cluster
 - better: use reconstruction error for cluster assignment
[\(Kambhatla and Leen, 1997; Hinton et al., 1997\)](#)
- These are limited by not having an overall density model
- [Tipping and Bishop](#): full probabilistic model using a mixture distribution in which components are probabilistic PCA
 - both discrete latent variables and continuous ones
 - [Kambhatla, N. and T. K. Leen](#) (1997). [Dimension reduction by local principal component analysis](#). *Neural Computation* **9**(7), 1493–1516.
 - [Hinton, G. E., P. Dayan, and M. Revow](#) (1997). [Modelling the manifolds of images of handwritten digits](#). *IEEE Transactions on Neural Networks* **8**(1), 65–74.
 - [Tipping, M. E. and C. M. Bishop](#) (1999a). [Mixtures of probabilistic principal component analyzers](#). *Neural Computation* **11**(2), 443–482.

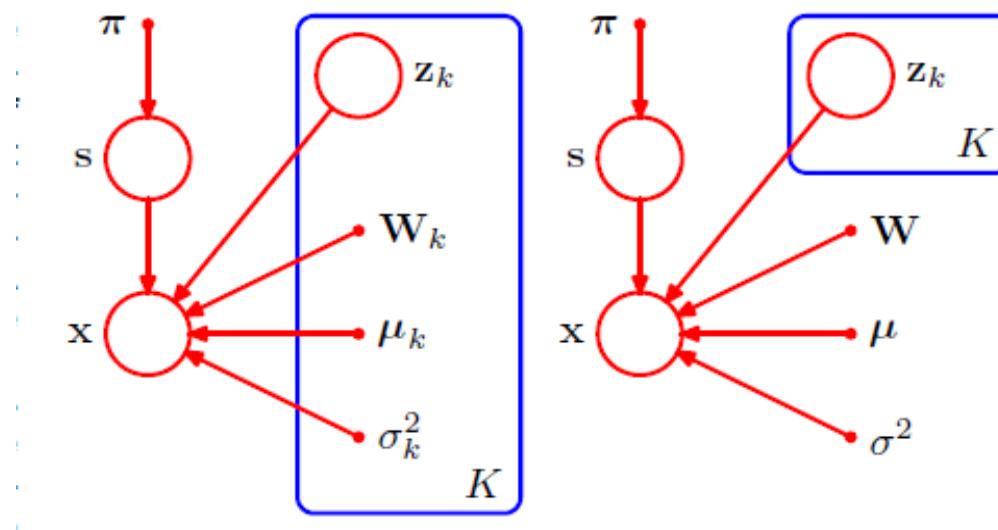
Mixture of PCA Models

- A fully Bayesian treatment, based on variational inference ([Bishop and Winn, 2000](#)): both the number of components in the mixture and the effective dimensionalities of the individual models are inferred from the data.
- There are many variants of this model in which parameters such as the W matrix or the noise variances are tied across components in the mixture, or in which the isotropic noise distributions are replaced by diagonal ones, giving rise to a [mixture of factor analyzers](#).
- Mixture of probabilistic PCA was extended hierarchically to produce an interactive data visualization algorithm.
 - [Bishop, C. M. and J. Winn](#) (2000). [Non-linear Bayesian image modelling](#). In *Proceedings Sixth European Conference on Computer Vision, Dublin*, Volume 1, pp. 3–17. Springer.
 - [Ghahramani, Z. and G. E. Hinton](#) (1996a). [The EM algorithm for mixtures of factor analyzers](#). Technical Report CRG-TR-96-1, University of Toronto.
 - [Ghahramani, Z. and M. J. Beal](#) (2000). [Variational inference for Bayesian mixtures of factor analyzers](#). In S. A. Solla, T. K. Leen, and K. R. Muller (Eds.), *Advances in Neural Information Processing Systems*, Volume 12, pp. 449–455. MIT Press.
 - [Bishop, C. M. and M. E. Tipping](#) (1998). [A hierarchical latent variable model for data visualization](#). *IEEE Transactions on Pattern Analysis and Machine Intelligence* **20**(3), 281–293.



Mixture of PCA Models

- Graphical models for the mixture of probabilistic PCA.
- On the right, all model components share all of the parameters (not very useful model)
- s is here the k –nomial latent variable that selects the mixture components – It has a prior governed by parameter π .



Principal Curves

- Alternative approach to use a single nonlinear model
- Principal Curves
 - Extension of PCA (that finds a linear subspace)
 - A curve is described by a vector-valued function $f(\lambda)$
 - Natural parametrization: the arc length along the curve
 - Given a point \hat{x} , we can find the closest point $\lambda = g_f(x)$ on the curve in terms of the Euclidean distance
 - A principal curve is a curve for which every point on the curve is the mean of all points in data space to project to it, so that
$$\mathbb{E}[x|g_f(x) = \lambda] = f(\lambda)$$
- There may be many principal curves for a continuous distribution
- Hastie et al: two-stage iterative procedure for finding principal curves
- Hastie, T. and W. Stuetzle (1989). Principal curves. *Journal of the American Statistical Association* **84** (106), 502–516.



Principal Curves

- For a given continuous density, there can be many principal curves.
- In practice, we are interested in finite data sets, and we also wish to restrict attention to smooth curves. [Hastie and Stuetzle \(1989\)](#) propose a two-stage iterative procedure for finding such principal curves, like the EM algorithm for PCA.
- The curve is initialized using the first principal component, and then the algorithm alternates between a data projection step and curve re-estimation step.
- **Projection step:** each data point is assigned to a value of λ corresponding to the closest point on the curve.
- **Re-estimation step:** each point on the curve is given by a weighted average of those points that project to nearby points on the curve, with points closest on the curve given the greatest weight.
 - [Hastie, T. and W. Stuetzle](#) (1989). [Principal curves](#). *Journal of the American Statistical Association* **84** (106), 502–516.



Principal Curves and Principal Surfaces

- When the subspace is constrained to be linear, the procedure converges to the first principal component
 - It is equivalent to the power method for finding the largest eigenvector of the covariance matrix.
- Principal curves can be generalized to multidimensional manifolds called principal surfaces
 - These have found limited use due to the difficulty of data smoothing in higher dimensions even for two-dimensional manifolds.



Multidimensional Scaling

- PCA is often used for the purpose of visualization
- Another technique with a similar aim: multidimensional scaling (MDS, [Cox and Cox 2000](#))
 - preserve as closely as possible the pairwise distances between data points.
 - involves [finding the eigenvectors of the distance matrix](#).
 - equivalent results to PCA when the distance is Euclidean.
 - can be extended to a wide variety of data types specified in terms of a [similarity matrix](#).

- Cox, T. F. and M. A. A. Cox (2000). [Multidimensional Scaling](#) (Second ed.). Chapman and Hall.
Machine Learning, University of Notre Dame, Notre Dame, IN, USA (Spring 2019, N. Zabaras)



Multidimensional Scaling

- Compute the matrix of squared pairwise similarities \mathbf{D} , $D_{ij} = \|\mathbf{x}_i - \mathbf{x}_j\|^2$.
- Compute $\mathbf{J} = \mathbf{I}_N - 1/N \mathbf{1}\mathbf{1}^T$ (where \mathbf{I}_N is the $N \times N$ identity function and N is the number of datapoints).
- Compute $\mathbf{B} = -\frac{1}{2} \mathbf{J} \mathbf{D} \mathbf{J}^T$.
- Find the L largest eigenvalues λ_i of \mathbf{B} , together with the corresponding eigenvectors \mathbf{e}_i .
- Put the eigenvalues into a diagonal matrix \mathbf{V} and set the eigenvectors to be columns of matrix \mathbf{P} .
- Compute the embedding as $\mathbf{X} = \mathbf{P} \mathbf{V}^{1/2}$.



Locally Linear Embedding

- Locally linear embedding (LLE, [Roweis and Saul 2000](#)).
- Compute the set of coefficients that best reconstruct each data point from its neighbors.
- Coefficients arranged to be invariant to rotation, translations, scaling

Characterize the local geometrical properties of the neighborhood.

- LLE maps the high-dimensional data to a lower dimensional subspace while preserving these coefficients.
- These weights are used to reconstruct the data points in low-dimensional space as in the high-dimensional space.
- Albeit non linear, LLE does not exhibit local minima.

- [Roweis, S. and L. Saul](#) (2000, December). [Nonlinear dimensionality reduction by locally linear embedding](#). *Science* **290**, 2323–2326.



Locally Linear Embedding

- Decide on the neighbours of each point (e.g., K nearest neighbours):
 - compute distances between every pair of points
 - find the k smallest distances
 - set $W_{ij} = 0$ for other points
 - for each point x_i :
 - ❖ create a list of its neighbours' locations \mathbf{z}_i ,
 - ❖ compute $\mathbf{z}_i = \mathbf{z}_i - \mathbf{x}_i$
- Compute the weights matrix W that minimizes $\varepsilon = \sum_{i=1}^N (x_i - \sum_{j=1}^N W_{ij} x_j)^2$ according to the constraints:
 - compute local covariance $\mathbf{C} = \mathbf{Z}\mathbf{Z}^T$, where \mathbf{Z} is the matrix of \mathbf{z}_i 's
 - solve $\mathbf{CW} = \mathbf{I}$ for W , where \mathbf{I} is the $N \times N$ identity matrix
 - set $W_{ij} = 0$ for non-neighbours
 - set other elements to $W / \sum W$
- Compute the lower dimensional vectors y_i that minimise $y_i = \sum_{i=1}^N (y_i -$



ISOMAP

- ❑ Isometric feature mapping (ISOMAP, [Tenenbaum et al. 2000](#))
- ❑ Goal: data projected to a lower-dimensional space using MDS

but dissimilarities defined in terms of the geodesic distances on the manifold

- ❑ Algorithm:
 - First defines the neighborhood using KNN
 - Construct a neighborhood graph with weights corresponding to the Euclidean distances
 - Geodesic distance approximated by the sum of Euclidean distances along the shortest path connecting two points
 - Apply MDS to the geodesic distances

- [Tenenbaum, J. B., V. de Silva, and J. C. Langford](#) (2000, December). [A global framework for nonlinear dimensionality reduction](#). *Science* **290**, 2319–2323.



Latent Trait Models

- We can also consider models having continuous latent variables together with discrete observed variables, giving rise to latent trait models ([Bartholomew, 1987](#)).
- In this case, the marginalization over the continuous latent variables, even for a linear relationship between latent and observed variables, cannot be performed analytically, and so more sophisticated techniques are required.
- [Tipping \(1999\)](#) uses variational inference in a model with a two-dimensional latent space, allowing a binary data set to be visualized analogously to the use of PCA to visualize continuous data.
 - Bartholomew, D. J. (1987). [Latent Variable Models and Factor Analysis](#). Charles Griffin.
 - [Tipping, M. E. \(1999\)](#). [Probabilistic visualisation of high-dimensional binary data](#). In M. S. Kearns, S. A. Solla, and D. A. Cohn (Eds.), *Advances in Neural Information Processing Systems*, Volume 11, pp. 592–598. MIT Press.



Latent Trait Models

- This model is the dual of the Bayesian logistic regression problem.
- In the case of logistic regression, we have N observations of the feature vector φ_n which are parametrized by a single parameter vector w , whereas in the latent space visualization model there is a single latent space variable x (analogous to φ) and N copies of the latent variable w_n .
- A generalization of probabilistic latent variable models to general exponential family distributions is described in [Collins et al. \(2002\)](#).
 - [Collins, M., S. Dasgupta, and R. E. Schapire](#) (2002). [A generalization of principal component analysis to the exponential family](#). In T. G. Dietterich, S. Becker, and Z. Ghahramani (Eds.), *Advances in Neural Information Processing Systems*, Volume 14, pp. 617–624. MIT Press.

Density Network

- An arbitrary distribution can be formed by taking a Gaussian random variable and performing a nonlinear transformation.
- This is exploited in a latent variable model called a density network ([MacKay, 1995](#); [MacKay and Gibbs, 1999](#)). The **nonlinear function is governed by a multilayered neural net**.
- If the network has enough hidden units, it can approximate a given nonlinear function to any desired accuracy.
- The marginalization over the latent variables to obtain the likelihood function, is no longer analytically tractable.
- The likelihood is approximated using Monte Carlo drawing samples from the Gaussian prior. The marginalization over the latent variables becomes a sum with one term for each sample.
- The procedure is costly: A large number of sample points may be required.

- [MacKay, D. J. C. \(1995\). Bayesian neural networks and density networks. *Nuclear Instruments and Methods in Physics Research, A* 354\(1\), 73–80.](#)
- [MacKay, D. J. C. and M. N. Gibbs](#) (1999). [Density networks](#). In J. W. Kay and D. M. Titterington (Eds.), *Statistics and Neural Networks: Advances at the Interface*, Chapter 5, pp. 129–145. Oxford University Press.

Generative Topographic Mapping



Generative Topographic Mapping

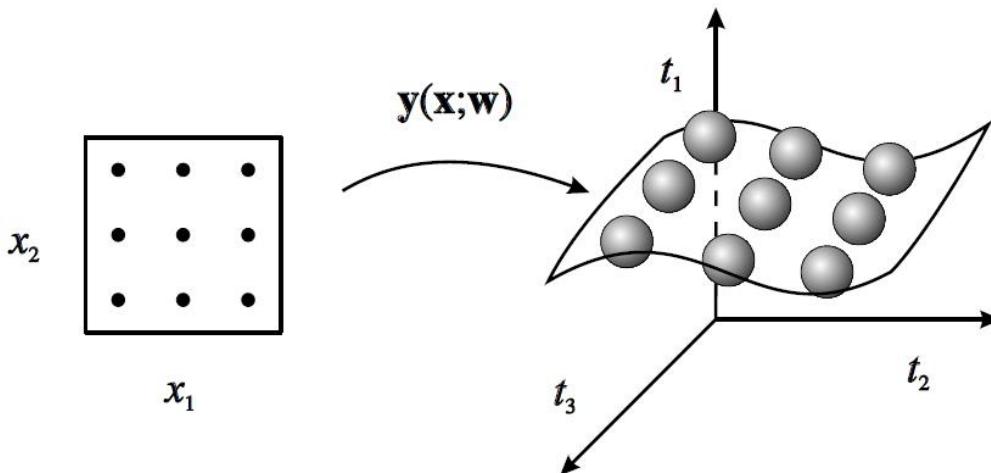
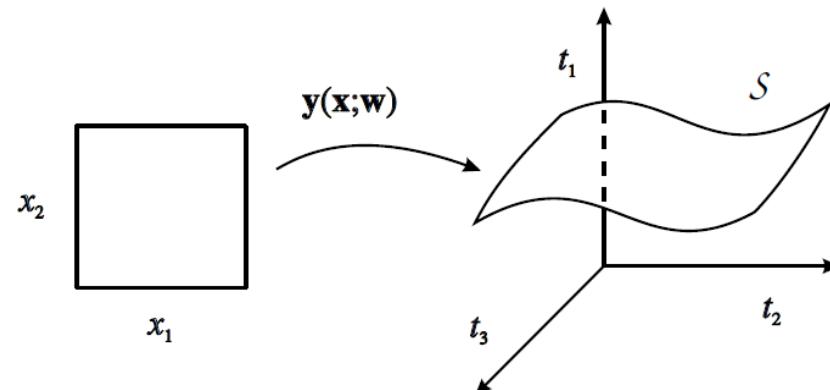
- If we consider more restricted forms for the nonlinear function, and make an appropriate choice of the latent variable distribution, then we can construct a latent variable model that is both nonlinear and efficient to train.
- The generative topographic mapping, or GTM ([Bishop et al., 1996](#); [Bishop et al., 1997](#); [Bishop et al., 1998](#)) uses a latent distribution that is defined by a finite regular grid of delta functions over the (typically two-dimensional) latent space.
- Marginalization over the latent space then simply involves summing over the contributions from each of the grid locations.

- [Bishop, C. M., M. Svensén, and C. K. I. Williams](#) (1997a). [GTM: a principled alternative to the Self-Organizing Map](#). In M. C. Mozer, M. I. Jordan, and T. Petche (Eds.), *Advances in Neural Information Processing Systems*, Volume 9, pp. 354–360. MIT Press.
- [Bishop, C. M., M. Svensen, and C. K. I. Williams](#) (1998b). [GTM: the Generative Topographic Mapping](#). *Neural Computation* **10**(1), 215–234.
- [Bishop, C. M., M. Svensen, and C. K. I. Williams](#) (1998a). [Developments of the Generative Topographic Mapping](#). *Neurocomputing* **21**, 203–224.



Generative Topographic Mapping

- The non-linear function $y(x; W)$ defines a manifold S embedded in data space given by the image of the latent-variable space under the mapping $x \rightarrow y$.

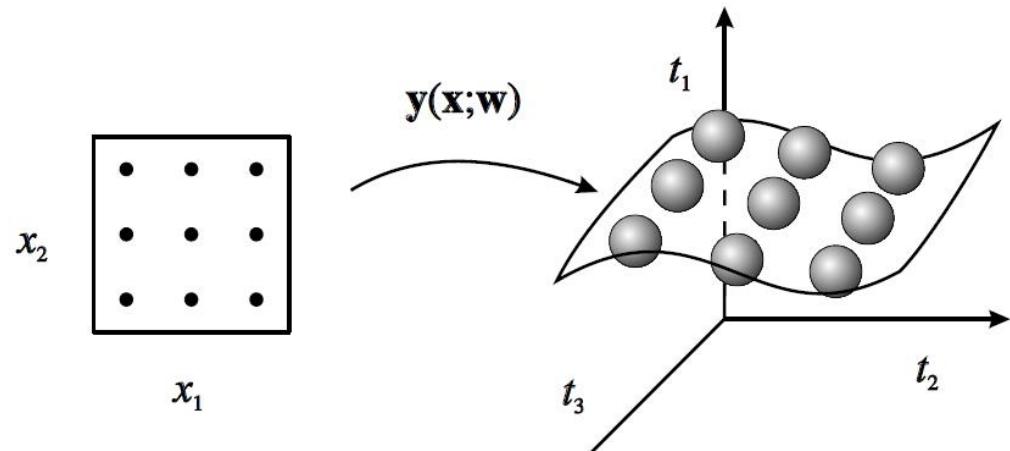


$$p(\mathbf{x}) = \frac{1}{K} \sum_{i=1}^K \delta(\mathbf{x} - \mathbf{x}_i)$$
$$p(\mathbf{t}|\mathbf{x}, \mathbf{W}, \beta) = \left(\frac{\beta}{2\pi} \right)^{D/2} \exp \left(-\frac{\beta}{2} \|\mathbf{y}(\mathbf{x}; \mathbf{w}) - \mathbf{t}\|^2 \right)$$

- We consider a prior distribution $p(\mathbf{x})$ consisting of a superposition of delta functions, located at the nodes of a regular grid in latent space. Each node \mathbf{x} is mapped to a corresponding point $y(\mathbf{x}; W)$ in data space, and forms the center of a corresponding Gaussian distribution.

Generative Topographic Mapping

- GTM results in a constrained Gaussian mixture model since the centers of the Gaussians cannot move independently but are constrained by $y(x; W)$.



$$p(\mathbf{x}) = \frac{1}{K} \sum_{i=1}^K \delta(\mathbf{x} - \mathbf{x}_i)$$

$$p(\mathbf{t}|\mathbf{x}, \mathbf{W}, \beta) = \left(\frac{\beta}{2\pi} \right)^{D/2} \exp \left(-\frac{\beta}{2} \|\mathbf{y}(\mathbf{x}; \mathbf{W}) - \mathbf{t}\|^2 \right)$$

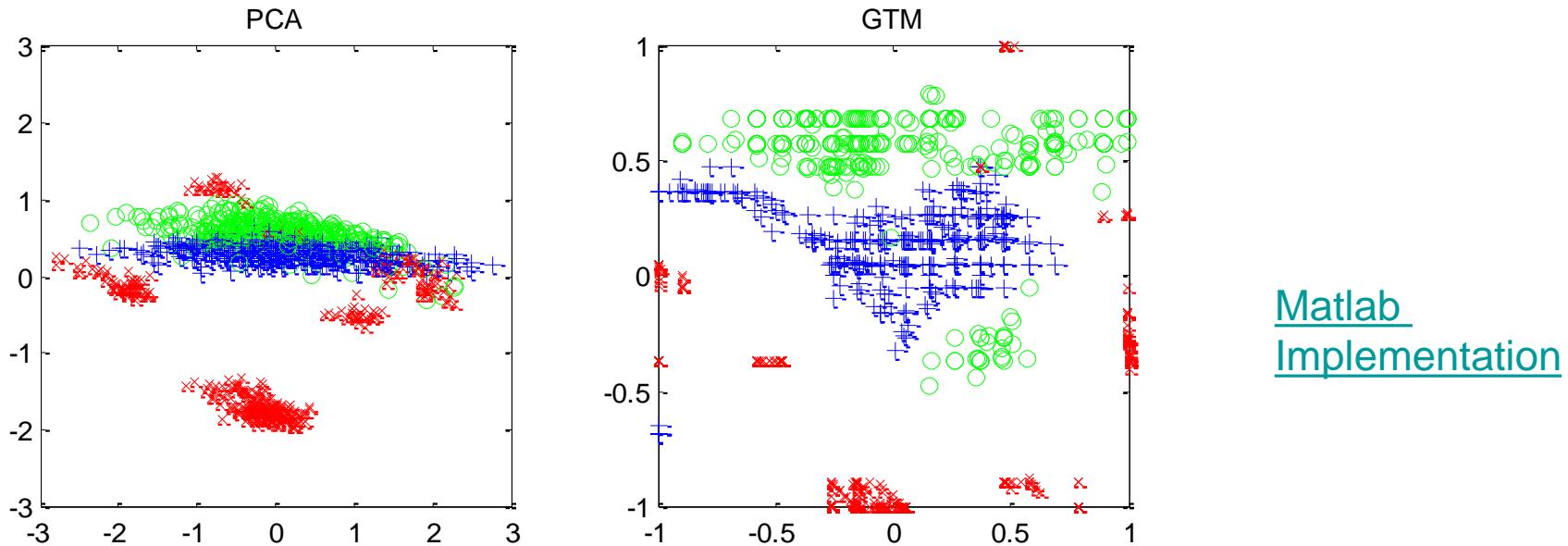
$$p(\mathbf{t}|\mathbf{W}, \beta) = \int p(\mathbf{t}|\mathbf{x}, \mathbf{W}, \beta)p(\mathbf{x})d\mathbf{x} = \frac{1}{K} \sum_{i=1}^K p(\mathbf{t}_i|\mathbf{x}_i, \mathbf{W}, \beta)$$

Generative Topographic Mapping

- The nonlinear mapping is given by a linear regression model that allows for general nonlinearity while being a linear function of the adaptive parameters.
- The limitation of linear regression models arising from the curse of dimensionality does not arise in the context of the GTM since **the manifold generally has 2 dimensions irrespective of the dimensionality of the data space.**
- As a consequence of these two choices, the likelihood function can be expressed analytically in closed form and can be optimized efficiently using the EM algorithm.
- The resulting GTM model **fits a 2D nonlinear manifold to the data set**, and by evaluating the posterior distribution over latent space for the data points, they can be projected back to the latent space for visualization purposes.

Generative Topographic Mapping

- The Figure shows a comparison of the oil data set visualized with linear PCA and with the nonlinear GTM.



- Each data point is plotted at the mean of its posterior distribution in latent space. The nonlinearity of the GTM allows clear separation between the groups of data points.

Self Organizing Maps (SOM)

- GTM can be seen as a probabilistic version of [the self organizing map](#), which also represents a 2D nonlinear manifold as a regular array of discrete points.
- SOM is reminiscent of K-means in that data are assigned to nearby prototype vectors that are subsequently updated.
- Initially, the prototypes are distributed at random, and during the training process they ‘self organize’ so as to approximate a smooth manifold.
- Unlike K-means, the SOM is not optimizing any well-defined cost function ([Erwin et al., 1992](#)) making it difficult to set the parameters of the model and to assess convergence.
- There is no guarantee that the ‘self-organization’ will take place as this is dependent on the choice of appropriate parameter values for any particular data set.

- [Kohonen, T.](#) (1982). [Self-organized formation of topologically correct feature maps](#). *Biological Cybernetics* **43**, 59–69.
- Kohonen, T. (1995). [Self-Organizing Maps](#). Springer.
- [Erwin, E., K. Obermayer, and K. Schulten](#) (1992). [Self-organizing maps: ordering, convergence properties and energy functions](#). *Biological Cybernetics* **67**, 47–55



Self Organizing Maps

- A set of K reference vectors \mathbf{z}_i is defined in the data space, in which each vector is associated with a node on a regular lattice in a 2D 'feature map' (analogous to the latent space of GTM).
- The algorithm begins by **initializing the reference vectors** (random values, equal to a random subset of the data points, or by using PCA). Each cycle of the algorithm then proceeds as follows.
- For every \mathbf{t}_n the corresponding 'winning node' $j(n)$ is identified, corresponding to \mathbf{z}_j having the smallest Euclidean distance $\|\mathbf{z}_j - \mathbf{t}_n\|^2$. The reference vectors are then updated by setting them equal to weighted averages of the data points given by ($h_{ij(n)}$ being a neighbourhood (Gaussian) function associated with node i)

$$\mathbf{z}_i = \frac{\sum_n h_{ij(n)} \mathbf{t}_n}{\sum_n h_{ij(n)}}$$

- The width of $h_{ij(n)}$ starts with a relatively large value and is gradually reduced after each iteration.



Extensions of the GTM Model

- By contrast, GTM optimizes the log likelihood function, and the resulting model defines a probability density in data space.
 - In fact, it corresponds to a **constrained mixture of Gaussians** in which the components share a **common variance**, and the means are constrained to lie on a smooth two-dimensional manifold.
 - This probabilistic foundation also makes it very straightforward to define generalizations of GTM ([Bishop et al., 1998](#)) such as
 - a **Bayesian treatment**,
 - dealing with **missing values**,
 - a principled **extension to discrete variables**,
 - use of Gaussian processes to define the manifold, or
 - a hierarchical GTM model ([Tino and Nabney, 2002](#)).
-
- [Tino, P. and I. T. Nabney \(2002\). Hierarchical GTM: constructing localized non-linear projection manifolds in a principled way. IEEE Transactions on Pattern Analysis and Machine Intelligence 24\(5\), 639–656.](#)
 - [Bishop, C. M., M. Svensen, and C. K. I. Williams \(1998a\). Developments of the Generative Topographic Mapping. Neurocomputing 21, 203– 224.](#)

Extensions of the GTM Model

- ❑ Because the manifold in GTM is defined as a continuous surface, not just at the prototype vectors as in the SOM, it is possible to compute the magnification factors corresponding to the local expansions and compressions of the manifold needed to fit the data set ([Bishop et al., 1997](#)) as well as the directional curvatures of the manifold ([Tino et al., 2001](#)).
- ❑ These can be visualized along with the projected data and provide additional insight into the model.

- [Tino, P. and I. T. Nabney \(2002\). Hierarchical GTM: constructing localized non-linear projection manifolds in a principled way. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **24**\(5\), 639–656.](#)
- [Bishop, C. M., M. Svensen, and C. K. I. Williams \(1998a\). Developments of the Generative Topographic Mapping. *Neurocomputing* **21**, 203–224.](#)